

PCTWORLD INTELLECTUAL PROPERTY ORGANIZATION
International Bureau**BI**

INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁶ : C12Q 1/68, C07H 21/04		A1	(11) International Publication Number: WO 99/43858
			(43) International Publication Date: 2 September 1999 (02.09.99)
(21) International Application Number: PCT/US99/04376		(74) Common Representative: McGINNIS, Robert, Owen; 1575 West Kagy Boulevard, Bozeman, MT 59715 (US).	
(22) International Filing Date: 26 February 1999 (26.02.99)			
(30) Priority Data: 60/076,102 26 February 1998 (26.02.98) US 60/076,182 27 February 1998 (27.02.98) US 60/086,947 27 May 1998 (27.05.98) US 60/107,673 7 November 1998 (07.11.98) US		(81) Designated States: AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, UA, UG, US, UZ, VN, YU, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).	
(63) Related by Continuation (CON) or Continuation-in-Part (CIP) to Earlier Applications US 60/076,182 (CIP) Filed on 27 February 1998 (27.02.98) US 60/086,947 (CIP) Filed on 27 May 1998 (27.05.98) US 60/076,102 (CIP) Filed on 26 February 1998 (26.02.98) US 60/107,673 (CIP) Filed on 7 November 1998 (07.11.98)		Published <i>With international search report.</i> <i>With amended claims and statement.</i>	
(71)(72) Applicants and Inventors: McGINNIS, Ralph, Evan [US/GB]; 27 Ladywell Prospect, Sawbridgeworth, Herts CM21 9PR (GB). McGINNIS, Robert, Owen [US/US]; 1575 West Kagy Boulevard, Bozeman, MT 59715 (US).			
(54) Title: TWO-DIMENSIONAL LINKAGE STUDY TECHNIQUES			
(57) Abstract <p>Versions of the invention are directed to methods (including software), apparatus, compositions of matter, and new uses of compositions of matter for a new type of association based linkage study technique using bi-allelic markers. In this new type of association based linkage study technique, the bi-allelic markers used in the new linkage studies are chosen so that the least common allele frequencies of the markers vary systematically over a range or subrange of least common allele frequency and the chromosomal location of the markers vary systematically over one or more chromosomes or chromosomal regions. And the bi-allelic markers are chosen so that the markers' chromosomal locations and least common allele frequencies vary systematically in an essentially independent manner. This selection of markers achieves a systematic distribution of the markers over a two-dimensional region having the orthogonal dimensions of chromosomal location and least common allele frequency. By using the two characteristics or two dimensions of marker chromosomal location and marker allele population frequency in this unique way, the power and systematic nature of genetic linkage studies using association based linkage tests is greatly increased. These unique two-dimensional linkage study techniques increase the power of association based linkage studies to localize trait causing genes or polymorphisms of modest effect such as human disease causing polymorphisms.</p>			

09/966,870 H 9

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

Two Dimensional Linkage Study Techniques

Technical Field

Versions of the present invention are in the field of molecular biology. some versions are specifically in the area of finding the chromosomal location of genes that cause genetic characteristics such as human disease.

Background

Introduction

Conventional linkage study techniques have limited power to localize trait causing genes (trait causing polymorphisms) of modest effect. such as many human disease polymorphisms. The two-dimensional linkage study techniques of this application are powerful new techniques for localizing genes (polymorphisms) especially of modest effect.

Chromosomes, heredity, genes, markers and alleles

Chromosomes are large molecules that carry the information for the inheritance of physical (genetic) characteristics or traits. In human beings for example, parents pass a copy of half of their chromosomes to their offspring during reproduction. By doing this, each parent passes some of his or her physical characteristics to his or her offspring. Any chromosome of a living creature is made of a large string-like molecule of DNA. Chromosomes are essentially very long strings of DNA. Genes are small pieces of a chromosome that cause or determine inherited genetic characteristics. (In this application, the term gene means a polymorphism that determines a genetic characteristic; the term does not mean an entire gene structure with a promoter region, introns, etc..) Markers are any segment of DNA on a chromosome which can be identified and whose chromosomal location is known (at least to some extent). Markers are like milestones along the very long string-like molecule of DNA which makes up a chromosome. Both a gene and a marker can come in different forms on different chromosomes. These different forms are known as different alleles and when a gene or marker comes in different forms it is said to be "polymorphic". For example, a bi-allelic marker comes in two (bi) different forms.

Linkage

If a gene allele and a marker allele occur as part of the genetic makeup of individuals more frequently than would be expected on the basis of chance, then it is possible to infer that the gene and the marker are linked. If a gene allele and a marker allele are inherited together more frequently than would be expected if the gene and the marker were on different chromosomes, then it is possible to infer that the gene and the marker are linked. Linkage of a gene and a marker usually occurs because the gene and the marker are close together on a chromosome. There are different degrees of linkage. Establishing linkage, especially strong linkage, between a gene and a marker can be very valuable. This is especially true if the precise location and other characteristics of the gene are not known. By establishing linkage, especially strong linkage, between a known marker and an unknown gene it is possible to locate the gene near to the chromosomal location of the known marker. This can be very

valuable if the gene is an important gene, such as a disease causing gene, and can help cure the disease.

Linkage Studies

Linkage studies are a method of establishing linkage between a marker and a gene or genes. Linkage studies are used to statistically correlate the occurrence of a genetic characteristic such as a disease (caused by a gene or genes) with a marker on a chromosome. One way this is done is by statistically correlating a specific allele of a marker with a genetic characteristic for a set of individuals by showing that individuals with the characteristic inherit the marker allele more often than individuals without the characteristic. The set of individuals is usually referred to as a sample of individuals. An example of a sample of individuals are people with a disease and similar people (matched controls) without the disease. Another example of a sample of individuals is a group of people, some of whom have the same disease; each of the people in the group being related to one or more of the other people in the group (i.e. families, sibships, pedigrees). The presence or absence of a marker allele in the chromosomal DNA of each individual is usually determined by genotype data at the marker for each individual.

There are different types of linkage study techniques, using different types of samples and different statistical measures of the correlation of a marker and a genetic characteristic. One example of a type of linkage study technique is the affected sib pair (ASP) test. Another example is the transmission disequilibrium test (TDT), which is an association based linkage test. This is a dynamic, changing area within the field of human genetics.

Linkage Studies and the "Scanning" of Chromosomal Regions

There are significant advantages in using several markers simultaneously to perform a linkage study with a genetic characteristic and a sample of individuals, especially when the relative positions of the markers on a chromosome are known. Such a linkage study allows searching for statistical evidence of linkage between markers in one or more regions of a chromosome or chromosomes and the gene or genes that determine the genetic characteristic. The results of the study for each marker can then be compared with the results for other markers, knowing the relative chromosomal positions of all the markers in the study. In this way, regions of a chromosome or even whole chromosomes can be "scanned" for evidence of linkage to a gene or genes causing a genetic characteristic. The relative positions of markers on chromosomes of a species of creatures is given by various kinds of chromosomal maps for the species. (There are several different kinds of marker maps, i.e. physical maps, genetic maps, radiation hybrid maps, etc.)

Sets of Markers for Linkage Studies and "Scanning" Chromosomes

An appropriate set of markers from a region of a chromosome can be chosen so that the region can be "scanned" for evidence of linkage of markers in the region to a gene or genes that cause a genetic characteristic. As explained above, this scanning is done by using the markers in linkage studies. Strong positive evidence for linkage of the markers (from the scanned chromosomal region) to a gene or genes responsible for a characteristic or trait is strong evidence that a trait-causing gene or genes is located within the chromosomal region.

Conventional Techniques for Choosing Sets of Markers to Scan Chromosomes with Linkage Studies

Conventional techniques choose sets of markers to scan a chromosomal region by choosing markers according to each marker's chromosomal location within the region. In a set of microsatellite markers described in 1994 for use in linkage studies, the markers were approximately evenly spaced, with average spacing between markers being 13 centiMorgans. The markers were distributed approximately evenly across the entire human genome (all human chromosomes) and were also selected because genotype data at the markers for individuals could be obtained by a semi-automated method.¹ A recent (1998) linkage study of the disease schizophrenia used a set of 310 microsatellite markers distributed approximately evenly across the entire human genome with average spacing of 11 centiMorgans between markers.² In a recent (1998) simulation of linkage studies to defend the practice of two-stage genome scanning, markers were spaced evenly every 10 cM (centimorgans) in an initial, sparser, first stage scan and evenly every 1 cM in a followup, denser, second stage scan.³ Following up positive linkage study results from chromosomal regions in a sparse, first stage scan with a second, denser scan that focuses on studying the regions with positive first-stage results is a common technique. In these conventional studies, as is common, markers were chosen to be about evenly spaced across the chromosomal regions studied. In this manner, as is conventional, a one dimensional structure such as an entire genome, a chromosome or a region of a chromosome is "covered" by markers in order to scan the entire genome, chromosome or chromosomal region with a linkage study. (These conventional techniques^{1,2,3} are not admitted to be prior art by their mention in this background.) (There is a possibly confusing, double meaning, of the term "marker map". It should be noted that a set of markers distributed along a chromosomal region, chromosome, or genome for linkage studies is also sometimes referred to as a "marker map" for use in chromosomal scanning by linkage studies. In addition, chromosomal or genetic maps of markers are also referred to as "marker maps".)

Conventional Techniques for Choosing Sets of Markers to Scan Chromosomal Regions are Essentially One Dimensional

Because DNA is a stringlike molecule, a chromosomal region(s), chromosome(s) and genome are essentially one dimensional in terms of the chromosomal location of markers and genes. As has been stated, conventional linkage study techniques scan a chromosomal region(s), chromosome(s) or genome by using markers distributed approximately evenly along the length of the chromosomal region(s), chromosome(s) or genome respectively. These conventional techniques focus primarily on the chromosomal location of markers used in a scan. These conventional techniques have an essentially one dimensional perspective.

Population Frequency of Marker Alleles and Gene Alleles

As described, chromosomal location of each marker is an important and unique characteristic of each marker and marker allele. Another characteristic of each polymorphic marker and each of the marker's

¹ Reed, et.al.: Chromosome-specific microsatellite sets for fluorescence-based, semi-automated genome mapping. Nature Genetics. July 1994: vol. 7: pp. 390-395.

² Levinson, et.al.: Genome Scan of Schizophrenia. Am J Psychiatry, June 1998: vol. 155: pp. 741-750.

³ Kruglyak, et. al.: Linkage Thresholds for Two-stage Genome Scans. Am J Hum Genet. 1998, vol. 62: pp. 994-996.

alleles is the population frequency of each marker allele. A population is a group (usually a large group) of individuals. A population frequency of a particular marker allele is the proportion of individual chromosomes in a population in which the particular marker occurs as the particular marker allele. For any bi-allelic marker, knowing the least common allele frequency of the marker establishes both of the allele frequencies of the marker. This is because the two allele frequencies of a bi-allelic marker sum to 1. Each gene allele also has a population allele frequency or allele frequency for short. Thus, each gene allele has a particular chromosomal location and allele frequency (for a particular population). In the case of an unknown gene, the gene's chromosomal location and allele frequencies are not specifically known.

Marker Allele Population Frequency in Conventional Linkage Study Scans

It is important to note that little attention was paid to the population allele frequencies of the markers used in the conventional linkage scans cited above. In the two studies cited above under conventional scanning techniques^{1,2}, marker allele frequency is referred to only peripherally as average marker heterozygosity, which is related to average marker allele frequency and the number of alleles (2, 3, 4, 5, etc.) at each marker. In the simulated scan cited above³, *the markers are stipulated to have four alleles that all have exactly the same allele frequency of 0.25 (heterozygosity 0.75). It is important to note that while the chromosomal location of the markers in all these conventional scans was systematically varied over the entire genome (all the human chromosomes), nothing was said about systematically varying the allele frequencies of the markers in any of the scans.* This is typical of conventional linkage study scans of genomes, chromosomes and chromosomal regions.

A Conventional View Of Bi-allelic Markers And Linkage Studies

We cite here a well known reference that discusses the conventional view of bi-allelic marker usefulness in linkage scans of chromosomes. In 1997 Kruglyak carried out computer simulations of the "information content" of markers that are part of various different marker maps.⁴ For bi-allelic markers his results showed that the optimum allele frequencies for bi-allelic markers used in linkage studies is 0.5/0.5 in order to achieve the greatest information content. However, allele frequency patterns other than the optimum 0.5/0.5 for bi-allelic markers gave acceptable levels of information content depending on the density of the marker map (or set of markers) chosen for the linkage study.

There are some important observations regarding this reference.⁴ First, *there is no advantage noted in this reference for choosing bi-allelic markers so that the set of chosen markers (or marker map) used for linkage studies is such that the markers systematically vary in allele frequency.* Thus, just as in the recent conventional linkage study scans cited above, there is no definite thought to using markers of systematically varying allele frequencies. The greatest information content is given by bi-allelic markers with allele frequencies close to the optimum of 0.5/0.5. Given the density of reasonably polymorphic SNPs predicted in this reference, at least one every 1 kb or 1,000 per cM, it is probable that even for quite dense maps, there will be so many acceptable SNPs available, that all of the SNPs in an appropriate marker map could have the optimum allele frequencies of approximately

⁴ Kruglyak: The use of a genetic map of biallelic markers in linkage studies. Nature Genetics. September 1997, vol.17, pp. 21-24.

0.5/0.5. Secondly, bi-allelic markers with lower least common allele frequencies, less than 0.3(0.7/0.3) or 0.2(0.8/0.2), are viewed unfavorably for linkage studies in this reference. Thirdly, the early version of the criterion of "information content" of markers used in this reference was based on sib pair analysis and the later, current version of the criterion, does not depend on any particular test for linkage.^{5, 6}

Thus, the criterion of information content in this reference, has never specifically employed the TDT (transmission disequilibrium test) or any association based test, whereas the two-dimensional linkage study techniques of this application are based on a completely different perspective of using association based tests. (This reference⁴ is not admitted to be prior art with respect to the present invention by it's mention in this background.)

Increased Power of the TDT (transmission disequilibrium test)

Characteristics of a new type of linkage test, the TDT (transmission disequilibrium test), were described in 1993. The inventor, R F McGinnis, was one of the authors of this reference.⁷ In 1996, Risch and Merikangas argued that conventional linkage analysis has limited power to detect genes of modest effect. And Risch and Merikangas attempted to illustrate the increased power of association based linkage tests such as the TDT over other types of conventional linkage tests.⁸ However, Risch and Merikangas' analysis was criticized by Muller-Myhsok and Abel as being based on the optimal assumption that the analyzed allele was the disease allele itself. Muller-Myhsok and Abel concluded that researchers should be aware that the power of association studies such as the TDT can be greatly diminished in more common, less optimal situations.⁹ In their response to Muller-Myshok and Abels' letter, Risch and Merikangas essentially agreed with the logic of Muller-Myshok and Abels' criticism. Risch and Merikangas stated that to a large extent, the expectation with respect to linkage disequilibrium across the genome is uncharted territory.¹⁰ (None of the references in this paragraph^{7,8, 9,10} is admitted to being prior art with respect to the present invention by their mention in this background.)

More Detailed Studies of the Power of the TDT

The inventor, R.E.McGinnis, has done extensive investigations on the power of the TDT. His observations and calculations of the increased power of the TDT in many situations have been

⁵ Kruglyak, et. al.: Complete Multipoint Sib-Pair Analysis of Qualitative and Quantitative Traits. Am J Hum Genet. 1995, vol. 57: pp. 439-454.

⁶ Kruglyak, et. al.: Parametric and Nonparametric Linkage Analysis: A Unified Multipoint Approach. Am J Hum Genet, 1996, vol. 58, pp. 1347- 1363.

⁷ Spielman, R.S., McGinnis, R.E., Ewens, W.J.: Transmission Test for Linkage Disequilibrium: The Insulin Gene Region and Insulin-dependent Diabetes Mellitus(1DDM). Am J Hum Genet. 1993. vol. 52. pp. 506-516.

⁸ Risch, N. and Merikangas, K.: The Future of Genetic Studies of Complex Human Diseases. Science. 13 September 1996, vol. 273, pp. 1516-1517.

⁹ Muller-Myshok, B. and Abel, L.: Technical Comments: The Future of Complex Diseases. Science. 28 February 1997, vol. 275, pp. 1328-1329.

¹⁰ Risch, N. and Merikangas, K.: Technical Comments: The Future of Complex Diseases. Science. 28 February 1997, vol. 275, p. 1330.

published.¹¹ In this paper a general framework for determining the power of the TDT in many different situations is presented. The analysis of Risch and Merikangas⁸ and others is shown by the inventor to be a special case of his general framework. His observations and calculations published in this paper have shown that the TDT has increased power in more common, less optimal situations as well as the less common, optimal situation cited by Muller-Myshok and Abel⁹. As opposed to the observation of Muller-Myhsok and Abel, the inventor's calculations indicate that association tests such as the TDT have increased power in typical situations even when the ratio m/p departs significantly from unity and, or the linkage disequilibrium between the analyzed (marker) allele and disease polymorphism is only half its maximum possible value. The inventor arrived at these conclusions independently and did not derive them from others.

A Major Conclusion Drawn by the Inventor about the TDT and Linkage Studies: Using Bi-allelic Markers of Systematically Varying Allele Frequencies Increases the Power of Linkage Studies Using the TDT

The inventor's calculations and observations about the increased power of the TDT in more common, less optimal situations led him to the conclusion that the power of linkage studies using the TDT is greatly increased under some conditions. Under some conditions, the power of the TDT in a linkage study using bi-allelic markers is greatly increased when each of one or more of the bi-allelic markers used in the study fulfill two criteria: (1) the allele frequencies of each of the one or more of the bi-allelic markers are similar (but not necessarily the same, or even approximately the same) as the allele frequencies of an unknown bi-allelic gene causing a disease under study; and (2) each of the one or more bi-allelic markers is in some degree of linkage disequilibrium with the gene. Thus for a typical linkage study using bi-allelic markers and the TDT, ***to increase the likelihood of conditions occurring that increase the power of the TDT in the linkage study, the bi-allelic markers used in the study are chosen so that the least common allele frequencies of the markers vary systematically over a range or subrange of least common allele frequency.*** This major conclusion of the inventor's research is quoted directly from his unpublished manuscript that was included with previously filed U.S. Provisional Patent Applications: "This example is typical and highlights perhaps the most important finding of this paper; namely the importance of using bi-allelic markers with heterozygosity similar to that of a bi-allelic disease gene. Indeed, since a majority of susceptibility loci may be bi-allelic, the judicious use of bi-allelic markers of both high, medium and low heterozygosity may be crucial in order to detect and replicate linkages to loci conferring modest disease risk." (page 25) (In this context the phrase "bi-allelic markers with heterozygosity similar to that of a bi-allelic disease gene" is essentially equivalent to "bi-allelic markers with individual allele frequencies similar to those of a bi-allelic disease gene" and "bi-allelic markers of both high, medium and low heterozygosity" is essentially equivalent to the phrase "bi-allelic markers whose least common individual allele frequencies are high, medium and low".)

¹¹ McGinnis, R.E.: Hidden Linkage: Comparison of the affected sib pair (ASP) test and transmission disequilibrium test (TDT). *Annals of Human Genetics*, 1998, vol. 62, pp. 159-179.

1 Systematically Varying Both Marker Chromosomal Location and Marker Allele Frequency of Markers in
2 Linkage Studies

3 The inventor's calculations and observations have demonstrated the increased power of the TDT in
4 more common, less optimal situations when a bi-allelic marker and bi-allelic gene have (1) similar but
5 not identical allele frequencies and (2) the marker and gene are in some degree of linkage
6 disequilibrium. Thus, for a typical linkage study using bi-allelic markers and the TDT, ***to increase the***
7 ***likelihood of both criteria (1) and (2) occurring for one or more markers, so as to increase the***
8 ***power of the TDT in the linkage study, the bi-allelic markers used in the study are chosen so***
9 ***that the least common allele frequencies of the markers vary systematically over a range or***
10 ***subrange of least common allele frequency AND the chromosomal location of the markers vary***
11 ***systematically over one or more chromosomes or chromosomal regions. And the bi-allelic***
12 ***markers are chosen so that the markers' chromosomal locations and least common allele***
13 ***frequencies vary systematically in an essentially independent manner.***

14 Two-dimensional Linkage Study Techniques

15 As has been stated, conventional linkage study scanning techniques use markers that are distributed
16 approximately evenly in the dimension of chromosomal location. These conventional, one dimensional,
17 scanning techniques focus primarily on the chromosomal location of markers used in a scan and give
18 little attention to the dimension of allele frequency.^{1, 2, 3}

19 One of the main implications of the inventor's work is to use a set of bi-allelic markers for a typical
20 linkage study using the TDT (or other association-based linkage test) wherein the chromosomal
21 locations and least common allele frequencies of the markers in the set systematically vary in an
22 essentially independent manner over the dimensions of chromosomal location and least common allele
23 frequency respectively. This is equivalent to using a set of bi-allelic markers for a linkage study scan
24 wherein the set of markers systematically scan or "cover" a two-dimensional region having dimensions
25 of chromosomal location and least common allele frequency. (Such a two-dimensional region can be
26 thought of as an area in an x-y plot or a group of squares on a chessboard.)

27 In addition, the inventor's calculations and observations indicate that bi-allelic markers having least
28 common allele frequencies less than 0.3, 0.2 or even less than 0.1 have an important place in linkage
29 studies using association based linkage tests. This is markedly different than Kruglyak's information
30 content evaluation of bi-allelic markers for use in linkage studies, in which bi-allelic markers with least
31 common allele frequencies less than 0.3 or 0.2 are viewed unfavorably.⁴

32 *In addition, the two-dimensional linkage study techniques do not necessarily favor using markers in a*
33 *scan that are about evenly spaced along a chromosome as in the conventional techniques. This is*
34 ***because conventional techniques suffer from a kind of one dimensional view or lack of depth***
35 ***perception. In the conventional techniques, a marker can look very close to a gene's location in***
36 ***terms of chromosomal location, but the marker can be very far from the gene's location in the***
37 ***new two dimensional view used by versions of the invention.***

38 ***It is as if the conventional 1D techniques look at a chessboard from on edge. Markers and a***
39 ***gene which are on different squares of the board, but in the same column of squares, look very***

close to each other when the board is looked at from on edge. But when the board is looked at from the top in 2D, two dimensions, markers which looked very close to each other and the gene before (when looking from on edge) can be seen to be very far from the gene.

Further Implications of the Two-dimensional Linkage Study Perspective

These two-dimensional techniques work when multiple genes cause a genetic characteristic and are effective in searching for these genes. A two-dimensional bi-allelic marker "covering" or scanning approach also increases the power of linkage studies using other association based linkage tests such as the AFBACmethod, the haplotype relative risk (HRR) method¹², and comparison of marker allele frequencies in disease cases and unrelated controls¹³. These references^{12, 13} are not admitted to being prior art with respect to the present invention by their mention in this background.)

Patents That May Be Helpful In Starting A Search Of The Background

Some patents that are in the same general areas as versions of the invention are cited here: US Patent Number 5,667,976 Solid supports for nucleic acid hybridization assays. Published International Application WO 98/20165 Biallelic Markers. Published International Application WO 98/07887 Methods for treating bipolar mood disorder associated with markers on chromosome 18 p. US Patent Number 5,552,270 Methods of DNA sequencing by hybridization based on optimizing concentration of matrix-bound oligonucleotide and device for carrying out same. No patent in this paragraph is admitted to being prior art with respect to the present invention by it's mention in this background.

¹² Falk CT and Rubenstein P: Haplotype relative risks: an easy reliable way to construct a proper control sample for risk calculations. Annals of Human Genetics, 1987, vol. 51, pp. 227-233.

¹³ Bell GI, Horita S and Karam JH: A polymorphic locus near the human insulin gene is associated with insulin-dependent diabetes mellitus. Diabetes, 1984, vol 33, pp. 176-183.

Two-Dimensional Linkage Study Techniques

Brief Description of Some Concepts Used By Versions of the Invention

Versions of the present invention make use of the novel concept of systematically covering a region on a two-dimensional map similar to an x-y graph with bi-allelic markers. The x axis on this map is the chromosomal location dimension and the y axis of the map is the least common allele frequency dimension. This two-dimensional map is called a CL-F map in this application. (CL stands for chromosomal location and F stands for least common allele frequency.) Each point on a CL-F map has two coordinates: a chromosomal location coordinate and a frequency coordinate. A point on a CL-F map is called a CL-F point.

Any one bi-allelic polymorphism (marker or gene) is viewed as being located at a particular CL-F point on a CL-F map. The chromosomal location of the polymorphism is the chromosomal location coordinate of the point. And the least common allele frequency of the polymorphism is the frequency coordinate of the point. The chromosomal location coordinate of a CL-F point is given in units of centiMorgans or base pairs or an equivalent thereof and the least common allele frequency coordinate of a CL-F point is given in units between 0 and 0.5 inclusive, such as 0.2.

Distances between any two CL-F points on a CL-F map are given in terms of two numbers: chromosomal location distance and frequency distance. The first number is the distance in the horizontal, chromosomal location direction. This first number is the chromosomal location distance. The second number is the distance in the vertical, frequency direction. This second number is the frequency distance. For example, the CL-F distance δ is given by two numbers δ_{CL} (chromosomal location distance) and δ_F (frequency distance). This is represented as $\delta = [\delta_{CL} \delta_F]$.

The "clustering" of bi-allelic markers near a particular CL-F point is discussed in terms of the number of markers within a particular CL-F distance of the point. For example, if each of N bi-allelic markers is separated from the point by a CL-F distance of less than or equal to δ , then the point is said to be N covered by the markers to within the distance δ . (N being an integer number.)

A region on a CL-F map is called a CL-F region. A CL-F region is a collection of one or more CL-F points. Some systematic methods of covering a CL-F region with bi-allelic markers are discussed in terms of the number of markers that are near each point in the region. For example, if each CL-F point in a CL-F region is N covered to within a CL-F distance δ by a subset of a set (or group) of bi-allelic markers, then the region is said to be N covered by the set (or group) of bi-allelic markers to within the distance δ .

A set (or group) of bi-allelic markers that cover a CL-F region or a CL-F point is referred to as a set (or group) of bi-allelic covering markers in this application.

The inventor discovered that when a bi-allelic marker and a bi-allelic gene are located close together on a CL-F map, then the power of association based linkage tests to detect linkage disequilibrium between the marker and a trait-causing gene (when present) increases greatly. Systematically covering a CL-F region that is the location of an unknown trait-causing bi-allelic gene with bi-allelic covering markers, therefore greatly increases the power of association based linkage tests to detect linkage disequilibrium (when present) between one or more of the covering markers and the gene.

1 A CL-F matrix is a matrix of rectangular cells of the same length and the same width on a CL-F map.

2 Stipulations that a certain number of covering markers are placed in each cell of the matrix is a method
3 of illustrating particular types systematic covering of a CL-F region with covering markers.

4 The evidence for linkage obtained from two-dimensional linkage studies is essentially two-dimensional
5 in nature and it is possible to use this two-dimensional information by essentially graphing quantitative
6 evidence for linkage as a function of position in the x-y plane. For example, if quantitative evidence for
7 linkage is represented in the z dimension of a typical three-dimensional x-y-z plot, wherein the x and y
8 dimensions are chromosomal location and least common allele frequency respectively, then it is
9 possible to conceptualize evidence for linkage as occurring in a "hump" or "humps" in the z dimension.

10 And it is possible to analyze the data to find the CL-F location (in the x-y plane) of the peak(s) of this
11 "hump(s)", thus helping to localize a trait causing gene to the CL-F locale of the peak(s) of the
12 "hump(s)".

13 Versions of the invention also make use of multi-allelic genes and/or markers. It is always possible to
14 combine the alleles of a multi-allelic polymorphism (marker or gene) so that the polymorphism acts
15 mathematically like it is a bi-allelic polymorphism. In effect, it is always possible to mathematically
16 transform a multi-allelic marker or gene to act bi-allelic. Similarly, two or more markers can always be
17 mathematically combined to form a mathematical marker that acts like a single bi-allelic marker. And
18 two or more genes can always be mathematically combined to form a mathematical gene that acts like
19 a single bi-allelic gene. In this application a mathematical bi-allelic marker formed mathematically from
20 one or more markers is called a bi-allelic marker equivalent or BME; and a mathematical bi-allelic gene
21 formed mathematically from one or more genes is called a bi-allelic gene equivalent or BGE.

22 The term true marker or gene is used to distinguish a marker or gene in the ordinary sense from a bi-
23 allelic marker equivalent (BME) or bi-allelic gene equivalent (BGE). The term true allele is used to
24 distinguish an allele in the ordinary sense from a mathematical allele of a BME or BGE. A mathematical
25 allele of a BME or BGE is referred to as an allele equivalent. An allele equivalent is a combination of
26 one or more true alleles or one or more haplotypes.

27 Versions of the invention make use of genes and/or markers, which are not exactly bi-allelic. These
28 genes or markers are approximately bi-allelic. A gene or marker that is approximately bi-allelic almost
29 always occurs in one of two allele forms, however, very rarely it occurs in a different allele form.

30 Various versions of the invention are for genotyping individuals at markers which systematically cover
31 CL-F regions or for obtaining sample allele frequency data (such as from pooled DNA) for a sample of
32 individuals for markers which systematically cover CL-F regions. Various versions of the invention are
33 for oligonucleotides used for genotyping individuals at markers which systematically cover CL-F regions
34 or are for obtaining sample allele frequency data (such as from pooled DNA) for a sample of individuals
35 for markers which systematically cover CL-F regions.

36 Definitions

38 For the purposes of the description and claims the terms used herein will have their generally accepted
39 definition unless otherwise specified.
40

11

1 The term **creature** means any organism that is living or was alive at one time. This includes both plants
2 and animals.

3 The term **species** is used in it's broadest sense and includes but is not limited to : 1)biological(genetic)
4 species,2) paleospecies (successional species), 3) taxonomic (morphological ; phenetic) species
5 including species hybrids such as mules, 4) microspecies (agamospecies) 5) biosystematic species(
6 coenospecies,ecosystem species)

7 A **genetic characteristic** is an observable or inferable inherited genetic characteristic or inherited
8 genetic trait including a biochemical or biophysical genetic trait, for example an inherited disease is a
9 genetic characteristic, a predisposition to an inherited disease is a genetic characteristic. A phenotypic
10 characteristic, phenotypic property or character is a genetic characteristic.

11 In this application, **the term gene** means a polymorphism that takes on one or more allele forms and
12 which causes or determines an inherited genetic characteristic or genetic trait. **The term gene** does not
13 mean an entire gene structure with a promoter region, a terminator region, introns, and other parts of an
14 entire gene structure. In this application the term gene means a polymorphism that determines or
15 causes an inherited genetic characteristic and that is part of an entire gene structure in some cases.
16 Each **genetic characteristic** of a creature is **determined** by one or more of the creature's **genes**,
17 wherein the term gene is defined as above.

18 A **segment** is a segment of a chromosome.

19 A **subrange** is a subrange of the least common allele frequency range 0 to 0.5 inclusive.

20 The **width** of a subrange is the difference between the upper and lower limits of the subrange. For
21 example, the width of the subrange 0.1 to 0.4 is $0.4 - 0.1 = 0.3$.

22 A **chromosomal location-least common allele frequency map** is a two-dimensional plot (similar to
23 an x-y graph) wherein the vertical axis(y axis) represents least common allele frequency and the
24 horizontal axis(x axis) represents chromosomal location. A chromosomal location-least common allele
25 frequency map is referred to as a **CL-F map**.

26 **Points on a CL-F map are referred to as CL-F points.** Points on a CL-F map have a chromosomal
27 location coordinate and a least common allele frequency coordinate. CL-F points represent possible
28 chromosomal location and least common allele frequency values for individual bi-allelic markers and
29 genes. Any particular point on a CL-F map is directly opposite a value on the map's least common
30 allele frequency axis(y axis) and is directly opposite a value on the map's chromosomal location axis(x
31 axis). These two values are the two coordinates of the particular point: (1) the chromosomal location
32 coordinate and (2) the least common allele frequency coordinate. A marker or gene located at a
33 particular point on a CL-F map is physically located at the chromosomal location given by the
34 chromosomal location coordinate of the point and the marker or gene's least common allele frequency
35 is the least common allele frequency coordinate of the point. These two coordinates are designated by
36 the term (x, y) wherein x is the value of the chromosomal location coordinate and y is the value of the
37 least common allele frequency coordinate.

38 A **particular CL-F map may be large or small.** For example it is possible for the chromosomal
39 location coordinates of CL-F points on a particular CL-F map to range over an entire chromosome (for
40 example human chromosome number 6). Alternatively it is possible for the chromosomal location

1 coordinates of CL-F points on a particular CL-F map to range over more than one chromosome, for
2 example all the human chromosomes, human chromosomes numbers 1 through 22 and X and Y.
3 Similarly it is possible for the chromosomal location coordinates of CL-F points on a particular CL-F
4 map to range over all the chromosomes of a species under study. Alternatively, it is possible for the
5 chromosomal location coordinates of CL-F points on a particular CL-F map to range over a very small
6 segment of chromosome, for example a segment of length 100,000 bp or less. Similarly it is possible for
7 the least common allele frequency coordinates of CL-F points on a particular CL-F map to range over
8 the entire least common allele frequency range 0 to 0.5. Alternatively it is possible for the least common
9 allele frequency coordinates of CL-F points on a particular CL-F map to range over a subrange or
10 subranges of the range 0 to 0.5, for example the subrange 0.1 to 0.2.

11 **If a bi-allelic polymorphism (marker or gene) is said to be located at a particular CL-F point then**
12 **the polymorphism's chromosomal location is the chromosomal location coordinate of the point and the**
13 **polymorphism's least common allele frequency is the least common allele frequency coordinate of the**
14 **point.**

15 **The chromosomal location distance between two CL-F points on a CL-F map is the absolute**
16 **difference between the two chromosomal location coordinates of the two points.**

17 **The frequency distance between two CL-F points on a CL-F map is the absolute difference between**
18 **the two least common allele frequency coordinates of the two points.**

19 **The CL-F distance between two CL-F points is given in terms of two parts or two components : (1)**
20 **chromosomal location distance and (2) frequency distance. This is denoted as $[D_{CL}, D_F]$, wherein D_{CL} is**
21 **the chromosomal location distance between the two points and D_F is the frequency distance between**
22 **the two points. For example $[500 \text{ bp}, 0.3]$ is an example of a CL-F distance.**

23 **If a first CL-F distance is less than or equal to a second CL-F distance then the chromosomal**
24 **location distance component of the first CL-F distance is less than or equal to the chromosomal location**
25 **distance component of the second CL-F distance AND the frequency distance component of the first**
26 **CL-F distance is less than or equal to the frequency distance component of the second CL-F distance.**
27 **For example if a first CL-F distance is $[x_1, y_1]$ and a second CL-F distance is $[x_2, y_2]$. And if the first CL-F**
28 **distance is said to be less than or equal to the second CL-F distance, then x_1 is less than or equal to x_2**
29 **AND y_1 is less than or equal to y_2 .**

30 **The term "bi-allelic covering marker(s)" or "covering marker(s)" is used to distinguish a particular**
31 **bi-allelic marker or particular bi-allelic markers from other markers. The term is being used simply to**
32 **avoid ambiguity. In general the term covering marker(s) can be thought of as a marker or markers**
33 **which have been chosen to cover or serve to cover a CL-F point or a CL-F region.**

34 **If a CL-F point is said to be N covered to within a CL-F distance δ by one or more bi-allelic**
35 **covering markers then the CL-F distance between each of N or more of the covering markers and the**
36 **point is less than or equal to δ . Wherein N is an integer number greater than or equal to 1.**

37 **If a CL-F point is said to be N covered to within a CL-F distance of about (or approximately) δ by**
38 **one or more bi-allelic covering markers then the CL-F distance between each of N or more of the**
39 **covering markers and the point is less than or equal to about (or approximately) δ . Wherein N is an**
40 **integer number greater than or equal to 1.**

1 **A CL-F region** is a group of CL-F points. A CL-F region is a region that is or can be represented on a
 2 CL-F map. A particular CL-F region may be large or small. For example the chromosomal location
 3 coordinates of CL-F points in a particular CL-F region can range over an entire chromosome (for
 4 example human chromosome number 6). Alternatively the chromosomal location coordinates of CL-F
 5 points in a particular CL-F region can range over more than one chromosome, for example all the
 6 human chromosomes, human chromosomes numbers 1 through 22 and X and Y. Similarly the
 7 chromosomal location coordinates of CL-F points in a particular CL-F region can range over all the
 8 chromosomes of a species under study. Alternatively, the chromosomal location coordinates of CL-F
 9 points in a particular CL-F region can range over only a small segment of chromosome, for example a
 10 segment of length 100,000 bp or less. Similarly the least common allele frequency coordinates of CL-F
 11 points in a particular CL-F region can range over the entire least common allele frequency range 0 to
 12 0.5. Alternatively the least common allele frequency coordinates of CL-F points in a particular CL-F
 13 region can range over only a very small subrange, for example the subrange 0.1 to 0.2 or less.

14 **The length of a CL-F region** is the largest chromosomal location distance between any two CL-F
 15 points in the region.

16 **The width of a CL-F region** is the largest frequency distance between any two CL-F points in the
 17 region.

18 **A CL-F region that is path connected** is contiguous and it is possible to draw a continuous path
 19 between any two points, wherein each point in the path is also in the region.

20 **If a CL-F region is said to be systematically covered by two or more bi-allelic covering markers**
 21 then each point in the region is within a small CL-F distance of one or more of the covering markers,
 22 wherein the magnitude of the small CL-F distance is such that there is increased power of an
 23 association based linkage test to detect evidence for linkage between one or more covering markers
 24 and a gene that is located at a point in the CL-F region, when linkage disequilibrium is present between
 25 the gene and one or more of the covering markers.

26 **If a CL-F region is said to be N covered to within a CL-F distance δ by one or more covering**
 27 **markers** then each point in the region is N covered to within the CL-F distance δ by the one or more
 28 covering markers. Wherein N is an integer greater than or equal to one.

29 **If a CL-F region is said to be N covered to within a CL-F distance of about (or approximately) δ**
 30 **by one or more covering markers** then each point in the region is N covered to within the CL-F
 31 distance of about (or approximately) δ by the one or more covering markers. Wherein N is an integer
 32 greater than or equal to one.

33 **The CL-F distance δ is known as the covering distance** if a CL-F point or CL-F region is N covered
 34 to within a CL-F distance δ .

35 **A CL-F covering distance δ has two components:** (1) a chromosomal location distance usually
 36 denoted by δ_{CL} and (2) a least common allele frequency distance (abbreviated as frequency distance)
 37 usually denoted by δ_F , i.e. $\delta = [\delta_{CL}, \delta_F]$.

38 **The length of a group of covering markers** is determined as follows. The absolute chromosomal
 39 location distance between each pair of markers in the group is determined. The greatest absolute

1 chromosomal location distance between each pair of markers in the group is the length of the group of
2 covering markers.

3 **A group of covering markers located on one chromosome can be ordered as a sequence of**
4 **markers** starting with the marker closest to one end of the chromosome and going toward the other
5 end of the chromosome. This is denoted for example as $m_1, m_2, m_3, \dots, m_{N-2}, m_{N-1}, m_N$, wherein N is
6 the number of markers in the group. (The chromosomal location distance between m_1 and m_N is greater
7 than the chromosomal location distance between any other pair of markers in the group and this
8 distance is the length of the group of markers.) **The chromosomal location distance between two**
9 **successive markers in the group**, i.e. between m_R and m_{R+1} , is a **chromosomal intermarker**
10 **distance**. (There are N-1 chromosomal intermarker distances for a group of N covering markers.)

11 **The average chromosomal intermarker distance** for a group is calculated by dividing the length of
12 the group by (N-1), wherein N is the number of covering markers in the group.

13 **The width of a CL-F region** is the largest frequency distance between any two CL-F points in the
14 region.

15 **The length of a CL-F region** is the largest chromosomal location distance between any two CL-F
16 points in the region.

17 **A segment-subrange pair** is the pair formed by pairing a segment of a chromosome and a subrange
18 of the least common allele frequency range 0 to 0.5.

19 The **term segment-subrange** is used as a short version of the term segment-subrange pair. (A
20 segment-subrange is a rectangular region on a CL-F map or a rectangular CL-F region, see below.)

21 **If one or more bi-allelic markers are said to be within(or in) a segment-subrange** then each of the
22 markers is located on (or in) the chromosomal segment of the segment-subrange(pair) and each of the
23 markers' least common allele frequencies is in the subrange of the segment-subrange(pair). (And each
24 of the markers is located within the rectangular region defined by the segment-subrange on a CL-F
25 map.)

26 Alternatively, if **a segment-subrange is said to contain one or more markers or to contain the**
27 **location of one or more markers** then each of the markers is located on (or in) the chromosomal
28 segment of the segment-subrange and each of the markers' least common allele frequencies is in the
29 subrange of the segment-subrange. (And each of the markers is located within or is within the
30 rectangular region on a CL-F map defined by the segment-subrange.)

31 **If one or more CL-F points are said to be within(or in) a segment-subrange** then each of the points
32 is located within the rectangular region defined by the segment-subrange on a CL-F map or on the
33 segment-subrange's borders.

34 **The length of a segment-subrange** is the length of the segment of the segment-subrange.

35 **The width of a segment-subrange** is the width of the subrange of the segment-subrange.

36 **The area of a segment-subrange** is the segment subrange's length multiplied by the segment
37 subrange's width.

38 **If a CL-F region is said to comprise a segment-subrange**, then each point in the segment-subrange
39 is in(or included in) the CL-F region.

1 If a **CL-F region** is said to comprise an area of greater than or equal to X multiplied by Y , then the
 2 CL-F region comprises one or more nonoverlapping segment-subranges, and the sum of the areas of
 3 the segment-subranges is greater than or equal to X multiplied by Y .

4 A **CL-F matrix** is a collection of segment-subranges, wherein each segment-subrange of the collection
 5 has the same width and the same length. Each segment-subrange in the collection (or the matrix) is a
 6 **CL-F matrix cell**. Any one CL-F matrix cell in a CL-F matrix shares two or more of the cell's borders
 7 with two or more other cells in the matrix. And all the cells in a CL-F matrix together form a single
 8 segment-subrange. A CL-F matrix is characterized by the length and the width of the cells in the
 9 denoted by length \times width, or $L_{MC} \times W_{MC}$, wherein L_{MC} is the length of each cell in the matrix and W_{MC} is
 10 the width of each cell in the matrix. A CL-F matrix is also characterized by the number of rows of cells,
 11 R_M , in the matrix. And a CL-F matrix is characterized by the number of columns of cells, C_M , in the
 12 matrix. There are two or more cells in a CL-F matrix. A CL-F matrix is also characterized by the point of
 13 origin of the matrix, denoted by (c_0, f_0) . The point of origin of a CL-F matrix is at any chromosomal
 14 location and c_0 takes on any reasonable value in an entire species genome. The point of origin of a
 15 CL-F matrix is at any one value in the least common allele frequency range 0 to 0.5. (A CL-F matrix is
 16 similar to the squares of a chessboard or to equal rectangular floor tiles that are all oriented in the same
 17 direction and cover a rectangular floor. One corner of the matrix is the matrix's point of origin.)
 18 The **width of each cell of a particular CL-F matrix** is any value greater than zero and less than 0.5.
 19 The width of a cell is often denoted by W_{MC} .

20 Any length in chromosomal location distance units is chosen for the **length of each cell of a particular**
 21 **CL-F matrix**. The length of a cell is often denoted by L_{MC} .

22 The **centerpoint** of a CL-F matrix cell is in the center of the cell. The centerpoints of a CL-F matrix form
 23 a **matrix centerpoint lattice**. Each point of a matrix centerpoint lattice is separated by a CL-F distance
 24 of $\{0, W_{MC}\}$ or $\{L_{MC}, 0\}$ from two or more neighboring centerpoints.

25 If **one or more bi-allelic markers are in(or within) the segment-subrange that is a CL-F matrix**
 26 **cell**, then each of the markers is in or within the CL-F matrix cell.

27 If **one or more CL-F points is in (or within) a CL-F matrix**, then each of the points is in or within a cell
 28 of the matrix.

29 If a **CL-F region comprises a CL-F matrix**, then each point that is in the matrix is also in the region.

30 If a **CL-F region is a CL-F matrix**, then the region consists of the points that are in the matrix.

31 If two CL-F matrix cells share a common border, then the **two CL-F matrix cells are in contact**.

32 If two CL-F matrix cells share a common corner, then the **two CL-F matrix cells are touching**. (Two
 33 cells that are in contact are also touching.)

34 If a **group of CL-F points is connected to within a CL-F distance $[X, Y]$** , then for any two points in
 35 the group, denoted p_1 and p_R , there is an ordered sequence of points in the group denoted $p_1, p_2,$
 36 $p_3, \dots, p_{R-2}, p_{R-1}, p_R$, R being an integer greater than or equal to 2, wherein the CL-F distance between
 37 each point in the sequence and the next point in the sequence is less than or equal to $[X, Y]$. The
 38 **distance $[X, Y]$ is the connecting distance**. (Put in simple terms if a group of points is connected to
 39 within $[X, Y]$, then there is a path between each pair of points in the group, the path consisting of a
 40 series of steps, wherein each step in the path is a movement between two points in the group that are

1 separated by a CL-F distance of less than or equal to $[X,Y]$. A simple group of points connected to
 2 within a CL-F distance of $[X,Y]$ is a group of three points, wherein each point in the group is within a CL-
 3 F distance of less than or equal to $[X,Y]$ of another point in the group. The concept of connectivity
 4 introduced here is similar to the basic concept of connectivity in mathematical graph theory.)
 5 **If a group of N markers is connected to within a CL-F distance $[X,Y]$, wherein N is an integer, then**
 6 **each of the markers is located at one point of group of N points, the group of N points being connected**
 7 **to within a CL-F distance $[X,Y]$.**
 8 **If two bi-allelic markers are said to be in extreme positive disequilibrium** then d is approximately
 9 equal to d_{\max} for the two markers, which for the purposes of this definition are designated marker M
 10 with least common allele A and marker m with least common allele B. Wherein according to standard
 11 usage, the disequilibrium coefficient (d) is defined by the equation $d=f(AB) - f(A)f(B)$ where $f(A)$ and $f(B)$
 12 are defined as the population frequencies of alleles A and B, respectively, and $f(AB)$ is the population
 13 frequency of the AB haplotype. And d_{\max} is defined as the maximum possible positive value of d
 14 assuming the allele frequencies of A and B are $f(A)$ and $f(B)$, and thus $d_{\max}= q-f(A)f(B)$ where q is the
 15 lesser of $f(A)$ and $f(B)$. (In this application d is used to represent the disequilibrium coefficient; the
 16 symbol δ is often used in scientific papers to represent the disequilibrium coefficient.)
 17 **If a pair of markers is said to be in extreme positive disequilibrium, then the two markers of the**
 18 **pair are in extreme positive disequilibrium.**
 19 **If a pair of bi-allelic markers is said to be redundant within distance D** then the two markers of the
 20 pair are in extreme positive disequilibrium and the two markers are located on the same chromosome
 21 and the two markers are located within a CL-F distance D of each other on a CL-F map, wherein D is a
 22 specified distance and D has two components, a chromosomal location distance component D_{CL} and a
 23 frequency distance component, D_F ; $D = [D_{CL}, D_F]$.
 24 **An allele equivalent (AE)** is a group of one or more "haplotype values" of one or more polymorphisms
 25 of the same type, either markers or genes. (For the purposes of this application a haplotype value of
 26 one polymorphism is equivalent to an allele value at the one polymorphism.)The group of haplotype
 27 values is then analyzed as if the group is a single allele at a bi-allelic polymorphism; the group of
 28 haplotype values acts as a single allele at a bi-allelic polymorphism; the collection of the one or more
 29 polymorphisms upon which the haplotype values are based acts as a bi-allelic polymorphism; the
 30 collection of one or more polymorphisms forms a bi-allelic **polymorphism equivalent (PE)** that acts as
 31 a bi-allelic polymorphism; **the polymorphism equivalent has(or possesses) the allele equivalent.**
 32 **The allele equivalent belongs to the polymorphism equivalent. In this application, each polymorphism**
 33 **equivalent is a bi-allelic marker equivalent(BME) or a bi-allelic gene equivalent(BGE).**
 34 **A bi-allelic marker equivalent (BME)** is one or more markers and a grouping of the haplotype values
 35 of the one or more markers into two groups (e.g. group I and group II)(For the purposes of this
 36 application a "haplotype value" of one marker is equivalent to an allele at the one marker). The one or
 37 more markers and the two groups of haplotype values of the one or more markers are then analyzed as
 38 if the one or more markers are a single bi-allelic marker with alleles I and II. Each group of the groups I
 39 and II is an allele equivalent. For example, a multi-allelic microsatellite marker has it's multiple alleles
 40 grouped into two groups and the microsatellite marker and these two groups of alleles then act

1 equivalent to a bi-allelic marker and are analyzed as if the microsatellite marker with the two groups is
2 bi-allelic (for an example of this see McGinnis, Ewens & Spielman, Genetic Epidemiology 1995 ; 12(6) :
3 637-40, which is incorporated herein by reference)
4 Also for example, two or more multi-allelic markers have their haplotypes separated into two groups of
5 haplotypes and the multi-allelic markers with their two groups of haplotypes are analyzed as if they
6 were a single bi-allelic marker.
7 For example bi-allelic marker A has alleles a and a* and bi-allelic marker B has alleles b and b*. Then
8 the four haplotypes ab, ab*, a*b* and a*b are grouped into two groups, for example group I: ab and a*b*
9 and group II : ab* and a*b. Then a BME formed by markers A and B takes on values of group I (or I) for
10 haplotypes ab or a*b* or group II (or II) for the haplotypes ab* or a*b ; and the two markers and the two
11 group values(I and II) are analyzed as though they form a single bi-allelic marker(the BME). The same
12 type of reasoning and procedure is extended to 3 or more bi-allelic markers, 3 or more bi-allelic marker
13 equivalents or 2 or more multi-allelic markers.
14 (Logically, of course, the genotype at a BME for an individual is determined by knowing the two
15 haplotype values at the one or more markers that form the BME for each of the individual's two
16 homologous chromosomes that carry the one or more markers. The genotype is then determined by
17 classifying each haplotype as belonging to group I or group II or the equivalent thereof. The three
18 possible genotype values at the BME are I / I, I / II, and II / II or the equivalent thereof.)
19 Similarly, a **bi-allelic gene equivalent (BGE)** is one or more genes and a grouping of all the haplotype
20 values of the one or more genes into two groups (e.g. group I and group II).
21 For the purposes of the description and claims, **the chromosomal location of a polymorphism**
22 **equivalent** is at any point on the smallest chromosomal segment that contains the one or more
23 polymorphisms that form the polymorphism equivalent(PE).
24 **The allele frequency of an allele equivalent (AE)** is determined as follows. An allele equivalent (AE)
25 is a group of haplotype values of one or more polymorphisms. The frequency of the allele equivalent is
26 the sum of the frequencies of the haplotype values in the group that makes up the allele equivalent.
27 For the purposes of the application, description, claims and definitions the term **true allele** is used to
28 distinguish an allele according to standard usage (i.e. at a single polymorphism) from an allele
29 equivalent (AE).
30 **The least common allele frequency of a bi-allelic polymorphism equivalent (BPE)** is determined
31 as follows. Each of the two groups(I and II) of the haplotype values of the one or more polymorphisms
32 which form the BPE is assigned a frequency. The frequency of I is the sum of the frequencies of the
33 haplotype values in group I. And the frequency of II is the sum of the frequencies of the haplotype
34 values in group II. The least of the frequency of I and the frequency of II is the least common allele
35 frequency of the BPE. If the frequency of I and the frequency of II are equal, then the least common
36 allele frequency of the BPE is the frequency of I or the frequency of II.
37 For the purposes of the description and claims, **the chromosomal location of a bi-allelic marker**
38 **equivalent (BME)** is at any point on the smallest chromosomal segment which contains the one or
39 more markers which form the BME.

1 **The chromosomal location distance from a BME to a CL-F point** on a CL-F map is the shortest
 2 chromosomal location distance from the CL-F point to any one of the one or more markers which form
 3 the BME.

4 **The least common allele frequency of a bi-allelic marker equivalent (BME)** is determined as
 5 follows. Each of the two groups(I and II) of the haplotype values of the one or more markers which form
 6 the BME is assigned a frequency. The frequency of I is the sum of the frequencies of the haplotype
 7 values in group I. And the frequency of II is the sum of the frequencies of the haplotype values in group
 8 II. The least of the frequency of I and the frequency of II is the least common allele frequency of the
 9 BME. If the frequency of I and the frequency of II are equal, then the least common allele frequency of
 10 the BME is the frequency of I or the frequency of II.

11 **The frequency distance from a BME to a CL-F point** on a CL-F map is the absolute difference
 12 between the least common allele frequency of the BME and the least common allele frequency
 13 coordinate of the CL-F point.

14 (If a CL-F point on a CL-F map is covered by one or more BMEs to within a distance δ , wherein $\delta = [\delta_{CL}$
 15 $, \delta_F]$, then the CL-F distance from each of the one or more BMEs to the CL-F point is less than or equal
 16 to δ . And the chromosomal location distance from one of the markers which form each BME to the CL-F
 17 point is less than or equal to δ_{CL} . And the frequency distance from each of the one or more BMEs to the
 18 CL-F point is less than or equal to δ_F .)

19 **A bi-allelic marker equivalent is in(or within) each CL-F matrix cell that contains the**
 20 **chromosomal location of the bi-allelic marker equivalent (BME).** (Since the chromosomal location
 21 of a bi-allelic marker equivalent (BME) is at any point on the smallest chromosomal segment which
 22 contains the one or more markers which form the BME, in some cases, a bi-allelic marker equivalent is
 23 in more than one CL-F matrix cell.)

24 For the purposes of the application, the term **true bi-allelic marker** is used to distinguish a bi-allelic
 25 marker with two alleles according to usual usage (i.e. at a single polymorphism) from a bi-allelic marker
 26 equivalent(BME). A true bi-allelic marker is not a bi-allelic marker equivalent (BME). The term **true bi-**
 27 **allelic polymorphism** is used to distinguish a bi-allelic polymorphism with two alleles according to
 28 usual usage from a bi-allelic polymorphism equivalent(BPE).

29 The term **true allele** of a true bi-allelic marker means an allele of a true bi-allelic marker.

30 **A polymorphism(marker or gene) which is exactly bi-allelic** has exactly two alleles and the sum of
 31 the frequency of each of the two alleles is 1; for example if the two alleles are A and B, then $f(A) + f(B)$
 32 $= 1$. A polymorphism that is exactly bi-allelic is a true bi-allelic polymorphism with exactly two true
 33 alleles or a bi-allelic polymorphism equivalent (BPE) with exactly two allele equivalents.

34 **A polymorphism(marker or gene) which is approximately bi-allelic** has three or more alleles. And
 35 the polymorphism has a first allele and a second allele; and the sum of the frequency of the first allele
 36 and the frequency of the second allele is approximately 1. And the frequency of the first allele and the
 37 frequency of the second allele is much greater than the sum of the allele frequencies of all the alleles of
 38 the polymorphism that are not the first or the second alleles. For the versions of the invention for bi-
 39 allelic polymorphisms (bi-allelic markers and bi-allelic genes) described herein, a polymorphism which
 40 is approximately bi-allelic is analyzed as if the polymorphism has only two alleles, the first allele and the

second allele. For the versions of the invention described herein, the least common allele frequency of a polymorphism which is approximately bi-allelic, is the least of the frequencies of the first and the second alleles of the polymorphism. A polymorphism which is approximately bi-allelic is a true polymorphism with true alleles (the allele frequencies of the true alleles conform to the stipulations of this definition) or is a bi-allelic polymorphism equivalent (BPE) with allele equivalents (the allele frequencies of the allele equivalents conform to the stipulations of this definition).

SNP stands for single nucleotide polymorphism.

A statistical linkage test based on allelic association is any mathematical test, mathematical computation or equivalent thereof which gives a quantitative estimate (or equivalent thereof) of evidence for linkage of a polymorphic marker and phenotypic trait (genetic characteristic) based on association between one or more of the alleles of the marker and the phenotypic trait in a sample of individuals of a population of a species. A statistical linkage test based on allelic association is any statistical test that detects or suggests linkage on the basis of allelic association. A statistical linkage test based on allelic association includes tests which suggest but do not prove linkage such as comparison of marker allele frequencies in disease cases and in unrelated controls. A statistical linkage test based on allelic association is also any test such as the TDT which may be regarded as "proving" linkage. (A statistical linkage test based on allelic association can, of course, give an estimate of the association of one or more allele equivalents of a marker equivalent and a genetic characteristic; see definition of BME above.) One aspect of a statistical linkage test based on allelic association is its potential use to calculate the probability, or equivalent thereof, that there is genuine association of one or more of the alleles of the marker and a genetic characteristic for the population as a whole (rather than just for the sample alone). A statistical linkage test based on allelic association is an association based linkage test. (The term **population** in this application is used in a statistical sense and means a group of individuals. The term **population** in this application is not used purely in the sense the term **population** is used in the field of population genetics.)

The term **sample** means a group of individuals which is a subset of a population.

In this application, **an allele is considered to be a piece of double stranded DNA** that is singular or distinctive for the allele. The piece of double stranded DNA that is distinctive for the allele contains the particular DNA sequence that distinguishes the allele from other alleles (alternate sequences) at the polymorphic site of interest plus two double stranded "flanking" DNA sequences, one flanking DNA sequence being on one side of the polymorphic site and the other flanking DNA sequence being on the other side of the polymorphic site.

Alternate strand of an allele: A double stranded piece of DNA that is distinctive for an allele consists of two pieces of single stranded DNA which are exactly complementary to one another. The two pieces of single stranded DNA are referred to as the two strands of the allele. Each of the two strands of the allele is the alternate of the other strand of the allele. For the purposes of this definition, the two strands are referred to as the first strand and the second strand. The alternate strand of the first strand is the second strand. And the alternate strand of the second strand is the first strand. Each strand of an allele is exactly complementary to the strand's alternate strand.

1 **An oligonucleotide** is either a single or double stranded oligonucleotide. The length of an
2 oligonucleotide ranges from a few bases or base pairs to approximately any number of bases or base
3 pairs in the DNA sequence of any allele.

4 **An oligonucleotide, either single or double stranded, is complementary** to an allele if the DNA
5 sequence of each strand of the oligonucleotide is exactly or approximately complementary to all or part
6 of the DNA sequence of one of the DNA strands of the allele and the oligonucleotide has utility in
7 identifying the allele by a hybridization reaction or equivalent thereof similar to as described below
8 under oligonucleotide technology.

9 **An allele is identified by a hybridization reaction with an oligonucleotide that is complementary**
10 **to the allele.** In this application there are two types of oligonucleotides that are complementary to
11 **an allele.** The two types of oligonucleotides complementary to an allele are identified as type(1) or
12 type(2).

13 A type (1) complementary oligonucleotide is complementary to the part of an allele's DNA sequence
14 that actually contains the allele's polymorphic site; and the type(1) complementary oligonucleotide has
15 utility to identify the allele by means of a hybridization reaction of the oligonucleotide to the part of the
16 allele's DNA sequence that actually contains the allele's polymorphic site. A hybridization reaction of a
17 type(1) oligonucleotide to the part of an allele's DNA sequence that actually contains the allele's
18 polymorphic site is a type (1) hybridization reaction.

19 A type (2) complementary oligonucleotide is complementary to an allele at a DNA sequence that flanks
20 (but does not contain) the allele's polymorphic site; and the type (2) complementary oligonucleotide has
21 utility to identify the allele by means of a hybridization reaction wherein the oligonucleotide hybridizes to
22 the allele at a DNA sequence that flanks (but does not contain) the allele's polymorphic site and
23 identification of the allele is subsequently achieved by extension of the oligonucleotide (and possibly
24 one or more other type(2)complementary oligonucleotides) across the polymorphic site with a DNA
25 polymerase such as occurs, for example, in a standard PCR (polymerase chain reaction). A
26 hybridization reaction of a type(2) oligonucleotide to an allele at a DNA sequence that flanks (but does
27 not contain) the allele's polymorphic site is a type (2) hybridization reaction.

28 Each version of **oligonucleotide technology** is a means to test for the presence (or absence) of each
29 of one or more true alleles of a group of true alleles in an individual's chromosomal DNA. The presence
30 or absence of any one true allele in the group is tested for by means of a type (1) or type (2)
31 hybridization reaction (or equivalent) with an oligonucleotide that is complementary(type(1) or type(2))
32 to the true allele. Put another way, the presence or absence of each true allele in the group is tested for
33 by means of a type(1) or type(2)hybridization reaction (or equivalent) with an oligonucleotide that is
34 complementary to each true allele in the group. There are many versions of oligonucleotide technology,
35 some of these versions are described in more detail below. (In this application, the term "chromosomal
36 DNA" includes chromosomal DNA obtained directly from an individual as well as DNA obtained as
37 amplification products using PCR and chromosomal DNA obtained directly from an individual.
38

1 **A physico-chemical signal** is any physical (including chemical) signal which is detected by human
2 senses or by apparatus. A physico-chemical signal includes, but is not limited to, (1) an electrical signal
3 such as is generated when oligonucleotides that are attached to a silicon chip hybridize with
4 complementary alleles, (2) a visual or optical signal such as is generated when oligonucleotides
5 attached to a glass slide hybridize with complementary alleles, (3) a signal (such as a dye color)
6 generated by the products of a PCR(polymerase chain reaction) such as when oligonucleotides that are
7 used as primers for PCR reactions hybridize with complementary alleles.

8 **The collection of true alleles of a group of one or more bi-allelic markers** is defined as consisting
9 of each true allele of each true marker in the group and each true allele of each haplotype that forms
10 each allele equivalent of each marker equivalent in the group.

11 **If a set of oligonucleotides is said to be complementary to a group of one or more bi-allelic**
12 **markers**, then each oligonucleotide in the set is type(1) or type(2) complementary to at least one of the
13 true alleles in the collection of true alleles of the group of one or more markers; and there is an
14 oligonucleotide in the set that is type(1) or type(2) complementary to each true allele in the collection of
15 true alleles of the group of one or more markers.

16 **Sample allele frequency data for a marker and a sample** is obtained by pooling DNA specimens
17 from individuals of the sample into one or more DNA pools. An allele frequency for each of the marker's
18 alleles is obtained for each DNA pool. In the case of a bi-allelic marker, determining the sample allele
19 frequency for one allele essentially determines the sample allele frequency for the other allele. (For
20 example, in some association based linkage studies, each DNA pool contains DNA from individuals of
21 the sample with the same or similar phenotype status.) (It is also possible to obtain sample allele
22 frequency for a marker and a sample by calculation using genotype data at the marker for each
23 individual in the sample.)

24 **Genotype data/sample allele frequency data** for a marker and a sample is (1)genotype data at the
25 marker for each individual of the sample, or (2)a combination of genotype data at the marker for one or
26 more individuals in the sample and sample allele frequency data for the marker for the sample, or
27 (3)sample allele frequency data for the marker for the sample. In the case of genotype data, DNA
28 specimens from individuals are tested individually to determine genotype. In the case of sample allele
29 frequency data DNA specimens from individuals are pooled, or sample allele frequency is calculated
30 using genotype data for each individual in the sample.

31 32 Description

33
34 For the versions of the invention described herein and the claims, **a bi-allelic genetic characteristic**
35 **gene or a bi-allelic gene** is a gene which is exactly bi-allelic or a gene which is approximately bi-allelic.
36 For the versions of the invention described herein and the claims, **a bi-allelic genetic characteristic**
37 **gene or a bi-allelic gene** is a gene which is a true bi-allelic gene or a bi-allelic gene equivalent (BGE).
38 A bi-allelic gene equivalent is exactly bi-allelic or approximately bi-allelic. A true bi-allelic gene is exactly
39 bi-allelic or approximately bi-allelic.

For the versions of the invention described herein and the claims, a **bi-allelic marker or a bi-allelic covering marker** is a marker which is exactly bi-allelic or a marker which is approximately bi-allelic. Each marker that is exactly bi-allelic is a true bi-allelic marker or a bi-allelic marker equivalent. And each marker that is approximately bi-allelic is a true bi-allelic marker or a bi-allelic marker equivalent (BME).

Process #1, A process for identifying one or more bi-allelic markers linked to a bi-allelic genetic characteristic gene in a species of creatures, comprising the steps of :

a)choosing two or more bi-allelic covering markers so that a CL-F region is systematically covered by the two or more covering markers;

b)choosing a statistical linkage test based on allelic association for each covering marker;

c)choosing a sample of individuals for each covering marker ;

d)obtaining genotype data/sample allele frequency data for each covering marker and the sample chosen for each covering marker, and obtaining phenotype status data for the genetic characteristic for each individual in the sample chosen for each covering marker;

e)calculating evidence for linkage between each covering marker and the gene using the statistical linkage test based on allelic association chosen for each covering marker and the genotype data/sample allele frequency data for each covering marker and using the phenotype status data for the genetic characteristic for each individual in the sample chosen for each covering marker obtained in d); and

f)identifying those covering markers as linked to the genetic characteristic gene which show evidence for linkage based on the calculations of step e.

The following is a more detailed description of process #1.

Process #1, A process for identifying one or more bi-allelic markers linked to a bi-allelic genetic characteristic gene in a species of creatures comprising the steps of :

a)choosing two or more bi-allelic covering markers so that a CL-F region is systematically covered by the two or more covering markers; Any method of systematically covering the CL-F region is acceptable. In this application, the systematic covering of a CL-F region in versions of the invention is described mathematically as the covering of a CL-F region, wherein the CL-F region is N covered to within a CL-F distance δ by two or more bi-allelic covering markers. For further details

1 regarding this step, see Detailed Description of the Systematic Covering of a CL-F Region Used In
2 Versions of the Invention below.

3
4 **b)choosing a statistical linkage test based on allelic association for each covering marker ;** The
5 statistical linkage test based on allelic association chosen for any one particular covering marker is any
6 statistical linkage test based on allelic association as defined in the definitions section. Statistical
7 linkage tests based on allelic association are described in the genetics and population genetics
8 literature and are known to those of ordinary skill in the art. Some examples of a statistical linkage test
9 based on allelic association are the TDT, Haplotype Relative Risk Method(HRR) and Allele Frequency
10 Comparison In Disease Cases Versus Unrelated Controls .It is possible for different statistical linkage
11 tests based on allelic association to be chosen for different covering markers. For purposes of technical
12 convenience, the same statistical linkage test based on allelic association is preferably chosen for each
13 covering marker.

14
15 **c)choosing a sample of individuals from the species for each covering marker ;** For the process
16 to be workable, the sample chosen for any one covering marker must be suitable for the statistical
17 linkage test of b) above chosen for the covering marker. Knowledge of a suitable sample for the
18 statistical linkage test chosen in b) above for the covering marker is within the understanding of a
19 person skilled in the art. For purposes of technical convenience, the same sample of individuals is
20 preferably chosen for each covering marker.

21
22 **d)obtaining genotype data/sample allele frequency data for each covering marker and the**
23 **sample chosen for each covering marker, and obtaining phenotype status data for the genetic**
24 **characteristic for each individual in the sample chosen for each covering marker;**
25 Sample allele frequency data for any one covering marker for the sample chosen for the covering
26 marker is obtained by pooling DNA from individuals of the sample into one or more DNA pools. It is also
27 possible to obtain sample allele frequency data for any one covering marker by calculation using
28 genotype data at the marker for each individual in the sample. Each DNA pool contains DNA from
29 individuals of the sample with the same or similar phenotype status. An allele frequency for each of the
30 marker's alleles is obtained for each pool. Genotype data/sample allele frequency data for any one
31 covering marker is (1)genotype data at the covering marker for each individual in the sample chosen for
32 the covering marker, or (2)a combination of genotype data at the covering marker for one or more
33 individuals in the sample chosen for the covering marker and sample allele frequency data for the
34 covering marker for the sample chosen for the covering marker, or (3)sample allele frequency data for
35 the covering marker for the sample chosen for the covering marker. The genotype data/sample allele
36 frequency data for any one covering marker must be suitable for the statistical linkage test based on
37 allelic association chosen for the covering marker in step b). It is possible to choose different types of
38 genotype data/sample allele frequency data for each covering marker. For purposes of technical
39 convenience, the same type of genotype data/sample allele frequency data (1), (2), or (3) is chosen for

each covering marker. Some examples of ways to practice this step is the use of technology cited under Oligonucleotide Technology (below) or mass spectrometry (such as MALDITOF).¹

e)calculating evidence for linkage between each covering marker and the gene using the statistical linkage test based on allelic association chosen for each covering marker and the genotype data/sample allele frequency data for each covering marker and using the phenotype status data for the genetic characteristic for each individual in the sample chosen for each covering marker obtained in d); and

f)identifying those covering markers as linked to the gene which show evidence for linkage based on the calculations of step e.

The meanings of steps d, e and f are within the understanding of those of ordinary skill in the art. Fine points of using a statistical linkage test based on allelic association as a measure of evidence for linkage are known to those in the art.¹¹

Process #1 described above is equivalent to localizing a genetic characteristic gene to a particular chromosomal location (i.e. a sub-region of a particular chromosome.) This is because markers which are linked to a gene are also physically close to the gene in terms of physical (chromosomal) location. To locate a gene causing the genetic characteristic of Process #1, the gene is localized to the approximate chromosomal location of one or more covering markers which are identified as showing evidence for linkage in step f).

Process#1A It is also possible to use Process #1 to localize a genetic characteristic gene to an approximate CL-F location(chromosomal location-least common allele frequency location). Such a process is expressed as follows:

Process#1A : A process for localizing a bi-allelic genetic characteristic gene in a species of creatures to a chromosomal location-least common allele frequency (CL-F) location, comprising the steps a), b), c), d) and e) of Process #1 and further comprising the step of:

f)localizing the gene to the chromosomal location-least common allele frequency (CL-F) location of one or more markers that show evidence for linkage based on the calculations of step e).

It is the teaching of this application that the strength of evidence for linkage increases as markers that are in linkage disequilibrium with a gene become close to the gene on a CL-F map. It is possible for step f) to be done by an individual plotting data by hand and examining the data. It is also possible for software to perform step f). It is possible for this step to include using the dependence of quantitative evidence for linkage of step e) on CL-F location. For example, if quantitative evidence for linkage calculated in step e) (of process #1 or #1A) is represented in the z dimension of a typical three-dimensional x-y-z plot, wherein the x and y dimensions are chromosomal location and least common allele frequency respectively, then it is possible to conceptualize evidence for linkage as occurring in a "hump" (or "humps") in the z dimension. And it is possible to use the evidence for linkage calculated in step e) of (process #1 or #1A) to find the CL-F location (in the x-y plane) of the peak(s) of a "hump(s)",

thus helping to localize a trait causing gene to the CL-F locale of the peak(s) of the "hump(s)". For example it is possible to use computer programming techniques that detect gradients such as, for example, linear or nonlinear programming techniques in mathematical optimization theory^{III} to find the peak(s) of a hump(s) in this step.

(Process #1A described above is equivalent to localizing a genetic characteristic gene to a particular chromosomal location (i.e. a sub-region of a particular chromosome.) This is because localizing a gene to a particular CL-F region also localizes the gene to a particular chromosomal region.)

Software

A computer program that executes each step of Process#1 is an example of Process#1. A computer program that executes each step of Process#1A is an example of Process#1A. A flowsheet illustrating programs that execute Process#1 and Process#1A is entitled Drawing #1(see drawing section). It is also possible for a computer program to execute any one of(or one or more combinations of) the steps of Process#1 or Process#1A. A person of ordinary skill in the art could write such a program without undue experimentation. The level of skill at computer programming in the art is great as evidenced by numerous computer programs. Some computer programs in the art are programs such as MAPMAKER/SIBS^{IV}, GENEHUNTER^V, LINKAGE^{VI}, and FASTLINK.^{VII}

Detailed Description of the Systematic Covering of a CL-F Region Used In Versions of the Invention

(see definitions section for meaning of CL-F region that is systematically covered). The CL-F region and covering markers are for a species and the one or more individuals are members of the species. The chromosomal location coordinate of each covering marker is based on information regarding the chromosomal location of each covering marker. One such source of information is chromosomal maps. Chromosomal maps are provided by such institutions as the Whitehead Institute or Marshfield Foundation for Biomedical Research. Chromosomal maps include, but are not limited to genetic maps, physical maps, and radiation hybrid maps.

The least common allele frequency coordinate of each covering marker is based on any reasonable information regarding the least common allele frequency of each covering marker. It is possible to use information from different populations for the allele frequencies of different covering markers. For example, it is possible for the least common allele frequencies of two different covering markers to be based on information from two different, but similar populations. For purposes of technical convenience, the least common allele frequency of each covering marker is based on information from the same population. One source of information on least common allele frequency is institutions which provide chromosomal maps such as the Whitehead Institute or Marshfield Foundation for Biomedical Research.

Systematic Covering Of A CL-F Region, Wherein A CL-F Region Is N Covered To Within A CL-F Distance δ By Two or more Bi-Allelic Covering Markers

In this application, the systematic covering of a CL-F region in versions of the invention is described mathematically as the covering of a CL-F region, wherein the CL-F region is N covered to within a CL-F distance δ by two or more bi-allelic covering markers. The covering markers are chosen so that the CL-F region is N covered to within the CL-F distance δ by using information regarding the chromosomal location and least common allele frequency of each covering marker.

It is possible for the chromosomal location component of δ to be as great as about any chromosomal length, computed by any method, for which linkage disequilibrium has been observed between any

polymorphisms in any population of the species. It is preferable in terms of increasing the power of a version of the invention for linkage studies that the chromosomal location component of δ be less than about the greatest chromosomal length, computed by any method, for which linkage disequilibrium has been observed between any polymorphisms in any population of the species. In general, the smaller the chromosomal location component of δ , the greater the power of a version of the invention for linkage studies.

It is possible for the frequency distance component of δ to be as great as about 0.2. (Depending on the penetrance ratio (r) or the disequilibrium between marker and gene, it is also possible for the frequency distance component of δ to be greater than 0.2 under some conditions as evidenced by Table 2 under Theory of Operation. So it is also possible for the frequency distance component of δ to be as great as about 0.25 or higher.) It is preferable in terms of increasing the power of a version of the invention for linkage studies that the frequency distance component of δ to be less than about 0.2. In general, the smaller the frequency distance component of δ , the greater the power of a version of the invention for linkage studies.

Linkage disequilibrium has been observed between polymorphisms separated by 10 to 12 cM in some homogeneous human populations. Therefore, it is possible for the chromosomal location distance component of δ to be as large as about 10 to 12 cM, about 10 to 12 million bp, or the equivalent thereof for homogeneous human populations. It is preferable in terms of increasing the power of a version of the invention for linkage studies in human populations that δ is less than or equal to about [1 million bp, 0.15] or the equivalent thereof. It is more preferable in terms of increasing the power of a version of the invention for linkage studies in human populations that δ is less than or equal to about [250,000 bp, 0.1] or the equivalent thereof.

In general, the smaller the magnitude of δ is in terms of either frequency distance, chromosomal location distance, or both, the greater the power of a version of the invention for linkage studies. In general, the greater N is, the greater the power of a version of the invention for linkage studies. Because the greater N is, the greater the chance that linkage is detected between one or more covering markers and a gene or genes. The largest that N is chosen is limited by the number of known markers in the neighborhood of the CL-F region and also by the distribution of the known markers. In general, the larger the CL-F region which is N covered, the greater the power of a version of the invention for linkage studies, because a larger region is scanned (covered). Less dense coverings wherein N is small, and the magnitude of δ is large also have technical and economic advantages for certain situations.

Specific types of CL-F regions that are N covered

Specific types of CL-F regions that are N covered are useful. For example, a rectangular CL-F region, a segment-subrange, that is N covered is used in an association based linkage study to test for the presence of a trait causing bi-allelic gene located within the segment-subrange. In the case in which a group of points is N covered to within a CL-F distance $[x,y]$ and the group of points is connected to within a CL-F distance of $[2x,2y]$ or less, then a path connected CL-F region is N covered to within the CL-F distance $[x,y]$.

1 A CL-F matrix is a device to illustrate and describe the systematic nature of special cases of CL-F
2 regions that are N covered. In the case in which there are N or more markers within each cell of a CL-F
3 matrix, then each point within the matrix is N covered to within the CL-F distance $[L_{CM}, W_{CM}]$, wherein
4 L_{CM} is the length of a matrix cell and W_{CM} is the width of a matrix cell. A choice of covering markers so
5 that approximately the same number of covering markers are in each cell of a CL-F matrix has utility in
6 that approximately the same amount of effort is expended on each subregion (cell) of the CL-F region
7 defined by the matrix in a linkage study using the covering markers. If the centerpoints of a CL-F matrix
8 (a matrix centerpoint lattice) are each N covered by a group of covering markers to within a CL-F
9 distance $[x,y]$, then each point in the matrix is N covered to within the CL-F distance $[2x,2y]$. A CL-F
10 matrix can be used as a device to help distinguish versions of the invention from prior art (to the extent
11 that there is prior art).

12 A requirement that the CL-F region that is N covered to within a certain CL-F distance comprise a
13 certain minimum area or segment-subrange with a certain minimum area is a special case of CL-F
14 regions that are N covered to within the certain CL-F distance. A requirement that the CL-F region that
15 is N covered to within a certain CL-F distance has a certain length or width is a special case of CL-F
16 regions that are N covered to within the certain CL-F distance. Each of these requirements is also a
17 device that can be used to help distinguish versions of the invention from prior art.

18 **A Note on the Equivalence of Working With Individual Alleles of Markers to Perform Two-**
19 **dimensional Linkage Studies and the CL-F approach using bi-allelic markers**

20 It is possible to conceptualize performing two-dimensional linkage studies wherein individual marker
21 alleles are used to cover a two-dimensional space, rather than individual bi-allelic markers. Any
22 individual marker allele is assigned a two-dimensional location consisting of the chromosomal location
23 of the marker and the allele frequency of the marker allele. Two-dimensional chromosomal location-
24 allele frequency spaces(or regions) are systematically covered by sets of covering alleles. Each
25 individual covering allele is tested for association with a genetic characteristic. Versions of inventions
26 using systematic chromosomal location-allele frequency(CL-AF) region coverings that are similar to
27 versions of the invention in this application are possible. Indeed these types of inventions have been
28 described in U.S.Provisional Patent Applications previously filed by the inventor.

29 However, such a conceptual framework and the resulting inventions are equivalent to the CL-F versions
30 approach used in this application. This is because any marker allele, A, that is used as a covering allele
31 can be made to be an allele equivalent of a bi-allelic marker equivalent(BME). So that a BME with allele
32 equivalents A and nonA is a bi-allelic marker with allele A. Therefore, any set of covering alleles that
33 systematically cover a two-dimensional CL-AF region is equivalent to a set of BMEs that systematically
34 cover an equivalent CL-F region. Testing each covering allele for association with a genetic
35 characteristic is exactly equivalent to testing each BME of a set of BMEs for evidence of linkage to a
36 gene using a statistical linkage test based on allelic association. Even testing for the presence or
37 absence of a covering allele in the chromosomal DNA of an individual is equivalent to genotyping the
38 individual at a BME. And determining a sample allele frequency for a covering allele, is equivalent to
39 determining the sample allele frequencies for a BME.

40

Example 1 of Process #1 is used for identifying markers linked to a disease gene.

Example 1 A process for identifying bi-allelic markers linked to a bi-allelic disease gene in human beings, comprising the steps of :

- a) choosing two or more bi-allelic covering markers so that a CL-F region is N covered to within a CL-F distance [250,000 bp, 0.1] or the equivalent thereof by the covering markers, wherein N is an integer number greater than or equal to 2 ;
- b) choosing the same statistical linkage test based on allelic association for each covering marker;
- c) choosing the same sample of individual human beings for each covering marker;
- d) obtaining genotype data at each covering marker for each individual in the sample and obtaining phenotype status data for the disease for each individual in the sample ;
- e) calculating evidence for linkage between each covering marker and the gene using the test chosen in step b) and the genotype data at each covering marker and the using the phenotype status data for the disease for each individual in the sample ; and
- f) identifying those covering markers as linked to the gene which show evidence for linkage based on the calculations of step e.

Apparatus Versions

General step by step descriptions of individual apparatus versions are given below.

Apparatus #1, an apparatus to practice process #1.

Apparatus #1, An apparatus for identifying bi-allelic markers linked to a bi-allelic genetic characteristic gene in a species of creatures, comprising :

- a) means for choosing two or more bi-allelic covering markers so that a CL-F region is systematically covered by the two or more covering markers;
- b) means for choosing a statistical linkage test based on allelic association for each covering marker;
- c) means for choosing a sample of individuals for each covering marker ;

d) means for obtaining genotype data/sample allele frequency data for each covering marker and the sample chosen for each covering marker, and for obtaining phenotype status data for the genetic characteristic for each individual in the sample chosen for each covering marker;

e) means for calculating evidence for linkage between each covering marker and the gene using the statistical linkage test based on allelic association chosen for each covering marker and the genotype data/sample allele frequency data for each covering marker and using the phenotype status data for the genetic characteristic for each individual in the sample chosen for each covering marker obtained in d); and

f) means for identifying those covering markers as linked to the gene which show evidence for linkage based on the calculations by means e).

More detailed description of Apparatus #1: Apparatus #1 is an apparatus to practice process #1. More details of the description of apparatus #1 are found under the description of Process #1 above. Any one of the means labeled a), b), c), d), e) or f) of apparatus #1 includes any means for automating or partially automating a step as step a), b), c), d), e) or f) respectively of process #1. An example of any one of the means in this paragraph labeled a), b), c), d), e), or f) is means comprising an appropriately programmed, suitable computer, the computer being supplied with proper data and instructions.

The means labeled d) of apparatus #1 for obtaining genotype data/ sample allele frequency data for each covering marker for the sample chosen for each covering marker includes any automated or partially automated means to obtain genotype data/ sample allele frequency data. An example of means to obtain genotype data/ sample allele frequency data is means using mass spectrometry.¹ Means to obtain genotype data/ sample allele frequency data that is automated or partially automated includes means comprising Oligonucleotide Technology described below.

Apparatus #1A, an apparatus to practice process #1A.

Apparatus#1A : An apparatus for localizing a bi-allelic genetic characteristic gene in a species of creatures to a chromosomal location-least common allele frequency (CL-F) region, comprising the means a), b), c), d) and e) of Apparatus #1 and further comprising the means of: f) means for localizing the gene to the approximate chromosomal location-least common allele frequency region (CL-F) of one or more markers that show evidence for linkage based on the calculations of means e).

An example of means f) is means comprising an appropriately programmed, suitable computer, the computer being supplied with proper data and instructions. Further details of this apparatus which practices process #1A are under process #1 and process #1A and Software(above).

Genotype data/Sample allele frequency data apparatus

An apparatus to obtain genotype data/sample allele frequency data similar to the data of the step d) of process #1 has great utility in that it is used to provide genotype data /sample allele frequency data for the more powerful two-dimensional linkage studies introduced in this application.

ApparatusGd/Safd#1: Genotype data/Sample allele frequency data apparatus: An apparatus for obtaining genotype data/sample allele frequency data for each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of a sample, comprising:

a) means for determining information on the presence or absence of each allele of each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of the sample, a CL-F region being systematically covered by the two or more bi-allelic covering markers; and

b) means for transforming the information of step a) into genotype data/sample allele frequency data for each marker of the group.

The CL-F region and covering markers are for a species and the one or more individuals are members of the species. Means for determining information on the presence or absence of each allele of each bi-allelic marker of the group in chromosomal DNA includes any means of determination. Means for determining information on the presence or absence of each allele of each bi-allelic marker of the group in chromosomal DNA includes means comprising oligonucleotide technology by using a set of oligonucleotides that is complementary to the group as discussed below. Information on the presence or absence of each allele in the chromosomal DNA is obtained using a DNA specimen from each of one or more individuals of the sample or by using one or more DNA pools of DNA specimens from two or more individuals of the sample. Any apparatus that obtains genotype data or sample allele frequency data (similar to the data of the step d) of process #1) by determining the presence or absence of each allele of each bi-allelic marker of the group in the chromosomal DNA of one or more individuals is an example of this version of the invention. Versions of this apparatus also obtain a combination of genotype data and sample allele frequency data similar to the data of the step d) of process #1. The details of step b) will be clear to those of ordinary skill in the art.

Each bi-allelic covering marker is a true bi-allelic or BME. Determining the presence or absence of each allele of each bi-allelic marker in the group includes determining the presence or absence of each allele equivalent of each bi-allelic marker equivalent(BME) in the group. Any method of systematically covering the CL-F region is acceptable. In this application, the systematic covering of a CL-F region in

versions of the invention is described mathematically as the covering of a CL-F region, wherein the CL-F region is N covered to within a CL-F distance δ by two or more bi-allelic covering markers. For further details regarding this, see Detailed Description of the Systematic Covering of a CL-F Region Used In Versions of the Invention above.

An example of ApparatusGd/Safd#1 Genotype data/Sample allele frequency data apparatus, a sample allele frequency apparatus:

Example 1 of ApparatusGd/Safd#1: An apparatus for obtaining genotype data/sample allele frequency data for each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of a sample, wherein the genotype data/sample allele frequency data is sample allele frequency data, comprising:

a) means for determining information on the presence or absence of each allele of each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA from one or more individuals of the sample, a CL-F region being N covered to within the CL-F distance [1.0 cM, 0.15] by the two or more bi-allelic covering markers, wherein N is an integer number greater than or equal to 1; and

b) means for transforming the information of step a) into sample allele frequency data for each marker of the group.

Example 2 of ApparatusGd/Safd#1: An apparatus for obtaining genotype data/sample allele frequency data for each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA of an individual, wherein the genotype data/sample allele frequency data is genotype data, comprising:

a) means for determining information on the presence or absence of each allele of each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA from an individual, a CL-F region being N covered to within the CL-F distance [12cM, 0.25] or the equivalent thereof by the two or more bi-allelic covering markers, wherein N is an integer number greater than or equal to 1; and

b) means for transforming the information of step a) into genotype data for each marker of the group.

(It should be noted that the following genotype apparatus is equivalent to Example 2 of ApparatusGd/Safd#1: Genotype Apparatus: An apparatus for genotyping an individual, comprising:

a) means to genotype an individual at two or more bi-allelic covering markers, a CL-F region being N covered to within the CL-F distance [12cM, 0.25] or the equivalent thereof by the two or more bi-allelic covering markers, wherein N is an integer number greater than or equal to 1.)

Genotyp data/Sample allele frequency data process

A process to obtain genotype data/sample allele frequency data similar to the data of the step d) of process #1 has great utility in that it is used to provide genotype data /sample allele frequency data for the more powerful two-dimensional linkage studies introduced in this application.

Description of the Genotype data/Sample allele frequency data process.

ProcessGd/Safd#1: Genotype data/Sample allele frequency data process: A process for obtaining genotype data/sample allele frequency data for each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of a sample, comprising:

a) determining information on the presence or absence of each allele of each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of the sample, a CL-F region being systematically covered by the two or more bi-allelic covering markers; and

b) transforming the information of step a) into genotype data/sample allele frequency data for each marker of the group.

The CL-F region and covering markers are for a species and the one or more individuals are members of the species. Determining information on the presence or absence of each allele of each bi-allelic marker of the group in chromosomal DNA includes any method of determination. Determining information on the presence or absence of each allele of each bi-allelic marker of the group in chromosomal DNA includes methods comprising oligonucleotide technology by using a set of oligonucleotides that is complementary to the group as discussed below. Information on the presence or absence of each allele in the chromosomal DNA is obtained using a DNA specimen from each of one or more individuals of the sample or by using one or more DNA pools of DNA specimens from two or more individuals of the sample. Any process that obtains genotype data or sample allele frequency data (similar to the data of the step d) of process #1) by determining the presence or absence of each allele of each bi-allelic marker of the group in the chromosomal DNA of one or more individuals is an example of this version of the invention. Versions of this process also obtain a combination of genotype data and sample allele frequency data similar to the data of the step d) of process #1. The details of step b) will be clear to those of ordinary skill in the art.

Each bi-allelic covering marker is a true bi-allelic or BME. Determining the presence or absence of each allele of each bi-allelic marker in the group includes determining the presence or absence of each allele equivalent of each bi-allelic marker equivalent(BME) in the group. Any method of systematically covering the CL-F region is acceptable. In this application, the systematic covering of a CL-F region in versions of the invention is described mathematically as the covering of a CL-F region, wherein the CL-F region is N covered to within a CL-F distance δ by two or more bi-allelic covering markers. For further

1 details regarding this, see Detailed Description of the Systematic Covering of a CL-F Region Used In
2 Versions of the Invention above.

3 **An example of ProcessGd/Safd#1Genotype data/Sample allele frequency data process, a**
4 **genotype data process:**

5 Example 1 of ProcessGd/Safd#1: A process for obtaining genotype data/sample allele frequency data
6 for each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA
7 of an individual, wherein the genotype data/sample allele frequency data is genotype data, comprising:

8 a) means for determining information on the presence or absence of each allele of each bi-allelic
9 marker of a group of two or more bi-allelic covering markers in the chromosomal DNA from an
10 individual, a CL-F region being N covered to within the CL-F distance [12cM, 0.25] or the equivalent
11 thereof by the two or more bi-allelic covering markers; wherein N is an integer number greater than or
12 equal to 1; and

13 b) means for transforming the information of step a) into genotype data for each marker of the group.

14 (It should be noted that the following genotype process is equivalent to Example 1 of
15 ProcessGd/Safd#1: Genotype Process: A process for genotyping an individual, comprising:

16 a)means to genotype an individual at two or more bi-allelic covering markers, a CL-F region being N
17 covered to within the CL-F distance [12cM, 0.25] or the equivalent thereof by the two or more bi-allelic
18 covering markers, wherein N is an integer number greater than or equal to 1.)

19
20 **Oligonucleotide technology**

21 Each version of oligonucleotide technology is a means to sense the presence or absence of each of
22 one or more true alleles of a group of true alleles in chromosomal DNA from one or more individuals by
23 means of a hybridization reaction with an oligonucleotide that is complementary to each of the one or
24 more true alleles (see definitions section). Thus versions of oligonucleotide technology are a means of
25 genotyping one or more individuals. And, versions of oligonucleotide technology are a means of
26 obtaining sample allele frequency data for one or more marker alleles for a sample of individuals using
27 pooled DNA from the individuals in the sample.

28 In Some Versions of Oligonucleotide Technology for Genotyping or Obtaining Sample Allele Frequency
29 Data, a Physico-chemical Signal is Generated when an Allele in Chromosomal DNA and a
30 Complementary Oligonucleotide Hybridize

31 Some versions of oligonucleotide technology for genotyping or for obtaining sample allele frequency
32 data use a sensor which includes one or more oligonucleotides which are complementary to an allele.
33 When the sensor is exposed to chromosomal DNA from an individual who carries the allele, the
34 oligonucleotides which are complementary to the allele hybridize with chromosomal DNA specimens of
35 the allele. The hybridization generates a physico-chemical signal which indicates the presence of the

allele in the chromosomal DNA of the individual. The lack of the physico-chemical signal indicates no (or negligible) hybridization and that the allele is not present in the chromosomal DNA of an individual.

Examples of oligonucleotide technology for genotyping, obtaining sample allele frequency data or genotype data/sample allele frequency data

Companies like Affymetrix are using high density arrays of oligonucleotides attached to silicon chips or glass slides to genotype DNA from one individual at thousands of bi-allelic markers.^{VIII} In some of these versions of oligonucleotide technology, the strength of hybridization of oligonucleotides that differ at only one base to DNA containing an SNP are compared to determine genotype.^{IX} Another version of oligonucleotide technology uses oligonucleotides as PCR (Polymerase Chain Reaction) primers to obtain genotype data.^X Other examples of oligonucleotide technology and its uses to obtain genetic information are included in the articles cited in the endnotes.^{XI} Versions of oligonucleotide technology obtain sample allele frequency data from pooled DNA or genotype data using oligonucleotides as PCR primers to obtain amplified reaction products that are detected by mass spectrometry. Another example of oligonucleotide technology is padlock probes.^{XII}

Other examples of oligonucleotide technology are minisequencing on DNA arrays, dynamic allele-specific hybridization, microplate array diagonal gel electrophoresis, pyrosequencing, oligonucleotide-specific ligation, the TaqMan system and immobilized padlock probes as presented at the First International Meeting on Single Nucleotide Polymorphism and Complex Genome Analysis.^{XIII}

Sets of Oligonucleotides for Genotyping at Bi-allelic Markers or Obtaining Sample Allele Frequency Data

A set of oligonucleotides that is complementary (see definitions) to a group of one or more bi-allelic markers has utility to determine genotype data at each of the markers in the group, including groups with BMEs and approximately bi-allelic markers.

Similarly, a set of oligonucleotides that is complementary to a group of bi-allelic markers has utility to obtain sample allele frequency data for each allele of each marker in the group.

In both cases, obtaining genotype data or sample allele frequency data, the same principle is used: a set of oligonucleotides that is complementary to a group of bi-allelic markers has utility to determine the presence or absence of each allele of each marker in the group in chromosomal DNA.

Using sets of oligonucleotides to obtain Genotype Data/Sample Allele Frequency Data for each marker of a group of bi-allelic markers, wherein the group of markers systematically cover a CL-F region

Genotype data/sample allele frequency data for each marker of a group of bi-allelic markers, wherein the group of bi-allelic markers systematically cover a CL-F region has great utility for use in the more powerful two-dimensional linkage studies introduced in this application. As described above under Oligonucleotide Technology, some sets of oligonucleotides have utility to determine genotype data at each bi-allelic marker of a group of one or more bi-allelic markers. Similarly, some sets of oligonucleotides have utility to obtain sample allele frequency data for each bi-allelic marker of a group of one or more bi-allelic markers. Therefore, the use of one or more copies of a set of oligonucleotides to obtain genotype data or sample allele frequency data for each bi-allelic marker of a group of one or

more bi-allelic covering markers, wherein the group of bi-allelic covering markers systematically cover a CL-F region has great utility.

A word to avoid confusion in terminology: in this application, a set of markers for use in genotyping is referred to as a set of oligonucleotides.

A set of oligonucleotides consisting of one or both strands of each allele of a group of one or more markers is a set of oligonucleotides that is complementary to the group of markers. (see definitions section) Such a set of oligonucleotides is in effect the group of markers themselves; and such a set of oligonucleotides has utility to determine genotype data at each marker in the group. So a group of markers (or set of markers) for use in obtaining genotype data or sample allele frequency data for each of the markers in the group is included in the descriptive phrase: "a set of oligonucleotides".

Description of Use set#1 D:

Use set#1 D The use of one or more copies of a set of oligonucleotides to determine genotype data/sample allele frequency data for each bi-allelic marker of a group of two or more bi-allelic covering markers for one or more individuals, wherein the group of covering markers systematically cover a CL-F region.

The CL-F region and covering markers are for a species and the one or more individuals are members of the species. An example of a set of oligonucleotides with utility to be used to determine genotype data/sample allele frequency data for each bi-allelic marker of a group of two or more bi-allelic covering markers is a set of oligonucleotides that is complementary to the group of markers. A set that is complementary to the group of markers is used to detect the presence or absence of each the alleles of the covering markers by means of a hybridization reaction as discussed under oligonucleotide technology. Thus a set that is complementary to the group of markers is used to determine genotype data/sample allele frequency data for each covering marker.

The use of one or more copies of a set of oligonucleotides to obtain genotype data or sample allele frequency data for each bi-allelic marker of a group of one or more bi-allelic covering markers, wherein the group of bi-allelic covering markers systematically cover a CL-F region are both examples of this version of the invention(Use Set#1D).

In this application, the systematic covering of a CL-F region in versions of the invention is described mathematically as the covering of a CL-F region, wherein the CL-F region is N covered to within a CL-F distance δ by two or more bi-allelic covering markers. For further details regarding this, see Detailed Description of the Systematic Covering of a CL-F Region Used In Versions of the Invention above.

Example 1S of Use set#1D: The use in genotyping one or more individuals, of one or more copies of a set of oligonucleotides, the set of oligonucleotides being complementary to a group of two or more bi-allelic covering markers, a CL-F region being N covered by the covering markers to within a CL-F distance of about [250,000 bp, 0.1] or the equivalent thereof, wherein N is an integer greater than or equal to two.

Composition of matter: Description of Comp set#1D:

Comp set#1D: One or more copies of a set of oligonucleotides, the set of oligonucleotides being

complementary to a group of two or more bi-allelic covering markers, wherein the group of covering markers systematically cover a CL-F region.

A set of oligonucleotides that is complementary to a group of two or more bi-allelic covering markers, wherein the group of covering markers systematically cover a CL-F region has great utility for use in the two-dimensional linkage study techniques introduced in this application. Such a set has utility in being used to genotype individuals or obtain sample allele frequency data or genotype data/sample allele frequency data as described above under Use set#1D. In this application, the systematic covering of a CL-F region in versions of the invention is described mathematically as the covering of a CL-F region, wherein the CL-F region is N covered to within a CL-F distance δ by two or more bi-allelic covering markers. For further details regarding this, see Detailed Description of the Systematic Covering of a CL-F Region Used In Versions of the Invention above.

Example 1Comp of Comp set#1D:

Example 1Comp: One or more copies of a set of oligonucleotides, the set of oligonucleotides being complementary to a group of two or more bi-allelic covering markers, a CL-F region being N covered by the covering markers to within a CL-F distance of about [1cM, 0.2] or the equivalent thereof, wherein N is an integer greater than or equal to one.

Redundancy of Covering Markers

Some versions of the invention make use of N coverings of CL-F regions by covering markers which limit (possibly to zero) the number of pairs of covering markers which are redundant within CL-F distance D, $D = [D_{CL}, D_F]$, wherein D is less than or equal to about δ , a CL-F covering distance. This limits the number covering markers which are separated by a CL-F distance of less than or equal to D (if the markers were placed on a CL-F map) which *will be in extreme positive disequilibrium with each other*. This limitation is done by requiring that less than or equal to R pairs of covering markers are redundant within distance D. Wherein R is an integer greater than or equal to 0 and less than or equal to about $N(N-1)/2$. When R is chosen to be zero, no pair of covering markers is redundant within distance D.

A preferable condition is that each bi-allelic covering marker within each small CL-F region (a small segment-subrange of length about δ_{CL} and width about δ_F the distance components of the covering distance δ) provides much new (i.e. non-redundant) information about linkage and association to any nearby bi-allelic gene. Under these conditions, testing each bi-allelic covering marker in each small CL-F region increases the likelihood of detecting linkage to a gene.

Limiting (including to zero) pairs of covering markers which are redundant within CL-F distance D (which is less than or equal to a covering distance δ) approaches and achieves this preferable condition. This limitation is not crucial to the functioning of a version of the invention, however, it has the advantage of reducing excess effort and increasing efficiency.

Polymorphism CL-F Display

Polymorphism CL-F display apparatus display the chromosomal location, least common allele frequency and identity of each polymorphism of one or more polymorphisms (markers or genes or both)

of one or more populations of one or more species on one or more two-dimensional graphs, each graph is similar to an x-y plot. The apparatus has utility including aiding in decisions regarding linkage studies and the interpretation of linkage study data.

The apparatus comprise means to display the chromosomal location, least common allele frequency and identity of each polymorphism of one or more polymorphisms (markers or genes or both) of one or more populations of one or more species on one or more two-dimensional graphs, each graph is similar to an x-y plot.

Each graph has two axes, one axis, the frequency axis, represents least common allele frequency and the alternate(or other) axis, the chromosomal location axis, represents chromosomal location. Each frequency axis of each graph is in units of population frequency. Each chromosomal location axis of each graph is in units of chromosomal location such as centimorgans, base pairs or the equivalent thereof.

The frequency axis of each graph spans the entire range 0 to 0.5 or a subrange of the range 0 to 0.5.

The chromosomal location axis of each graph spans the chromosomal locations on one or more segments of one or more chromosomes of a species, each of the one or more segments is a size from the equivalent of a base pair in length to the length of an entire chromosome (or the equivalent thereof).

Each point on each graph is directly opposite a value on the frequency axis of each graph. The value on the frequency axis directly opposite each point on each graph is the frequency coordinate of each point on each graph. Each point on each graph is directly opposite a value on the chromosomal location axis of each graph. The value on the chromosomal location axis directly opposite each point on each graph is the chromosomal location coordinate of each point on each graph.

Each graph displays the chromosomal location and least common allele frequency of each polymorphism of one or more polymorphisms. Each polymorphism displayed on each graph is assigned a graph location on each graph.

The graph location of each polymorphism displayed on each graph is typical of the use of x-y plots. The graph location assigned to each polymorphism on each graph is a point. The chromosomal location coordinate of the point assigned as the graph location to any one polymorphism is equal (or approximately equal) to the chromosomal location of the polymorphism. And the frequency coordinate of the point assigned as the graph location to any one polymorphism is equal (or approximately equal) to the least common allele frequency of the polymorphism.

The apparatus comprise means for displaying one or more two-dimensional graphs. Each graph comprises, the identity and graph location of one or more polymorphisms assigned a location on each graph. And the apparatus comprise means for displaying one or more graphs wherein the viewer chooses the species, population, polymorphisms, span of the frequency axis and span of the chromosomal location axis of the one or more graphs ; in versions of the apparatus, the means of this sentence comprises a computer.

The apparatus comprise means for storing and updating data on the chromosomal location and least common allele frequency of one or more polymorphisms of one or more populations of one or more species and means for storing chromosomal location and least common allele frequency data on newly discovered polymorphisms.

Versions of the apparatus comprise means for printing each of the one or more graphs.

Theory of Operation / Best Mode

Systematically Varying Both Marker Chromosomal Location and Marker Allele Frequency of Markers in Linkage Studies

The inventor's calculations and observations have demonstrated the increased power of the TDT in more common, less optimal situations when a bi-allelic marker and bi-allelic gene have (1) similar but not identical allele frequencies and (2) the marker and gene are in some degree of linkage disequilibrium. Thus, for a typical linkage study using bi-allelic markers and an association based linkage test, ***to increase the likelihood of both criteria (1) and (2) occurring for one or more markers, so as to increase the power of an association based linkage test in a linkage study, the bi-allelic markers used in the study are chosen so that the least common allele frequencies of the markers vary systematically over a range or subrange of least common allele frequency AND the chromosomal location of the markers vary systematically over one or more chromosomes or chromosomal regions. And the bi-allelic markers are chosen so that the markers' chromosomal locations and least common allele frequencies vary systematically in an essentially independent manner.***

(In the Theory of Operation/ Best Mode Section the traditional symbol used in scientific papers for the disequilibrium coefficient, δ , is used. This should not be confused with the symbol δ used for the covering distance in the remainder of the application. The symbol d is used for the disequilibrium coefficient in the sections of the application other than the Theory of Operation/Best Mode Section.)

The theory of operation is based on the mathematical observation that the TDT and other association-based tests for linkage are increased in power as the frequencies of the disease-causing allele of a bi-allelic gene and the positively associated allele of a linked bi-allelic marker become similar in magnitude. The inventor made this observation as a result of deriving the equation shown below for P_t (this is Equation 2 in the unpublished manuscript submitted for publication in December 1996 and in

published paper by RE McGinnis in the Annals of Human Genetics vol 62, pp. 159-179, 1998).

$$P_t = .5 + (1 - 2\theta) \left[\frac{c_1 c_4 - c_2 c_3}{H} \right] \left\{ p^2 \left(\frac{\alpha^2 - \beta^2}{4} \right) + 2p(1-p) \left(\frac{(\alpha + \beta)^2 - (\beta + \gamma)^2}{16} \right) + (1-p)^2 \left(\frac{\beta^2 - \gamma^2}{4} \right) \right\}$$

Equation 2

P_t may be regarded as the size of the "signal" which is given by the TDT to indicate that a tested marker is linked to a disease-causing gene. The more P_t is elevated above 0.5 (baseline), the greater is the evidence for linkage or "power" provided by the association-based linkage test known as the TDT.

Table 2 in the unpublished manuscript filed with previous US Provisional Patent Applications (see below) illustrates how signal strength increases substantially as the frequencies of disease-causing allele and positively associated marker allele become similar in magnitude. As noted on pages 24 and 25 of the unpublished manuscript (see below), Table 2 assumes that the frequency (p)

1 of the disease-causing allele is fixed at $p=.1$ while the frequency (m) of the positively associated marker
 2 allele varies ($m=.5, .3, .2, .1, .05$). Note that when the level of disequilibrium (or association) between
 3 the bi-allelic marker and bi-allelic disease gene is fixed (in this case either $\delta=\delta_{\max}$ or $\delta=\frac{1}{2}\delta_{\max}$), the
 4 signal strength of P_t progressively increases as m decreases from $m=.5$ to $m=.1$ (the same frequency
 5 as the disease allele, i.e., $p=.1$). For example, in the section of Table 2 for $r=5$, note that when $\delta=\frac{1}{2}$
 6 δ_{\max} , P_t is .548 at $m=.5$ and then steadily increases to .572 ($m=.3$), .597 ($m=.2$), .648 ($m=.1$) and then
 7 starts to decrease again as m departs from $m=p=.1$ (i.e. $P_t=.636$ at $m=.05$). As noted on pages 24-25
 8 (below) of the unpublished manuscript, the TDT chi-square statistic (assuming a sample size of 200
 9 families) is such that the signal strength at $m=.5$ ($P_t=.548$) does not produce a statistically significant
 10 evidence for linkage ($p\text{-value} > 0.5$) while the doubling of signal strength at $m=.2$ ($P_t=.597$) produces
 11 very strong statistical evidence for linkage by the TDT ($p\text{-value} < 0.005$). This sort of substantial
 12 increase in power is also true of other association-based linkage tests as the frequencies of the
 13 disease-causing allele and associated marker allele become more similar in magnitude.
 14

1 Table 2(Footnotes for Table 2 are on next page)

2 Effect of penetrance ratio (r), disequilibrium (δ) and marker heterozygosity (m) on magnitude
3 of P_t and P_s

			Magnitude of P_t			Magnitude of P_s		
			δ_{\max}^a	$\frac{1}{2} \delta_{\max}^b$	$\delta=0$	δ_{\max}^a	$\frac{1}{2} \delta_{\max}^b$	$\delta=0$
6	r=2	m=.5	.526	.513	.500	.505	.505	.504
7		m=.3	.541	.521	.500	.508	.506	.504
8		m=.2	.558	.531	.500	.511	.508	.504
9		m=.1	.595	.555	.500	.518	.512	.504
10		m=.05	.589	.552	.500	.517	.511	.504
11								
12	r=5	m=.5	.596	.548	.500	.543	.540	.539
13		m=.3	.633	.572	.500	.561	.548	.539
14		m=.2	.666	.597	.500	.575	.556	.539
15		m=.1	.719	.648	.500	.600	.573	.539
16		m=.05	.696	.636	.500	.589	.571	.539
17								
18	r=10	m=.5	.656	.577	.500	.595	.587	.584
19		m=.3	.702	.612	.500	.623	.600	.584
20		m=.2	.736	.644	.500	.644	.612	.584
21		m=.1	.785	.703	.500	.673	.635	.584
22		m=.05	.750	.684	.500	.652	.628	.584
23								
24	r= ∞	m=.5	.740	.617	.500	.700	.680	.673
25		m=.3	.791	.663	.500	.743	.700	.673
26		m=.2	.826	.703	.500	.772	.716	.673
27		m=.1	.870	.770	.500	.807	.744	.673
28		m=.05	.816	.741	.500	.763	.730	.673
29								

Footnotes for Table 2

a,b Value of δ that is maximal (δ_{\max}) and half-maximal ($\frac{1}{2} \delta_{\max}$), as determined by the heterozygosity of the marker (m) and disease locus ($p=.1$)

Importance of disequilibrium and marker heterozygosity (i.e. marker allele frequency) in detecting linkage

When the heterozygosity (i.e. allele frequencies) of a bi-allelic marker and bi-allelic disease locus are fixed, ($P_S = .5$) and $|P_t - .5|$ are both maximized at the most positive or most negative possible value of δ (δ_{\max} , δ_{\min}), as demonstrated in the published paper. This maximization of χ^2_{asp} and χ^2_{ldt} is intimately connected to M_S and M_t (defined in equations 1 and 2) since: (a) these are the only two factors in P_S and P_t that are influenced by δ and (b) M_S and $|M_t|$ are maximal and equal to each other when δ is extreme (δ_{\max} or δ_{\min}). Furthermore, as explained in the published paper, M_S is a measure of the proportion of informative (A/B) parents who are also informative (D/d) at the disease locus. Therefore, maximizing M_S (and, by implication, $|M_t|$) is equivalent to *minimizing* the proportion of A/B parents who are homozygous (D/D or d/d) at the disease locus. Such homozygous D/D or d/d parents contribute evidence *against* linkage since they transmit marker alleles A and B to affected offspring with equal probability; thus, minimizing their proportion among A/B parents being tested for linkage has the effect of maximizing χ^2_{asp} and χ^2_{ldt} .

Nevertheless, when bi-allelic markers have a specific (i.e. fixed) heterozygosity different from that of a bi-allelic disease locus, some A/B parents must be homozygous at the disease locus, even when δ is extreme. But if marker heterozygosity is variable, the proportion of A/B parents who are D/D or d/d approaches *zero* as marker heterozygosity approaches that of the disease locus and as δ approaches δ_{\max} or δ_{\min} . Consequently, the most extreme values of P_t and P_S , and highest values of χ^2_{ldt} and χ^2_{asp} are found when marker and disease locus have equal heterozygosity and $\delta = \delta_{\max}$ or $\delta = \delta_{\min}$.

Example illustrating the importance of marker heterozygosity (i.e. allele frequency)

To illustrate the importance of marker heterozygosity and disequilibrium, Table 2 shows P_t and P_s values when the frequency (p) of disease allele D is constant at 0.1, but the frequency (m) of marker allele A varies between $m=.5$ (maximum marker heterozygosity) and $m=.1$ (equal heterozygosity at marker and disease loci). The table assumes mode of inheritance is additive, and separate sections of the table show the results when the penetrance ratio (r) is 2, 5, 10 or ∞ . For each value of r , an individual line in the table represents constant marker heterozygosity ($m=.5$, $.3$, $.2$, or $.1$) and from left-to-right on each line, one sees P_t and P_s values when $\delta=\delta_{\max}$, $\delta=\frac{1}{2}\delta_{\max}$, and $\delta=0$, the value of δ_{\max} being determined by the particular values of m and p [$\delta_{\max}=p(1-m)$]. As noted in Appendix I of the published paper, when $p<m$ and $p<(1-m)$, as in this example, the most extreme values of P_t and P_s must occur at $\delta=\delta_{\max}$. This can be seen in each line of the table by the steady increase in both P_t and P_s as one moves from $\delta=0$ to $\delta=\delta_{\max}$, with every line also showing $P_t > P_s$ at $\delta=\delta_{\max}$ and most lines showing $P_t > P_s$ at $\delta=\frac{1}{2}\delta_{\max}$.

Most remarkable, however, are the sizeable increases in P_s and even greater increases in P_t as marker heterozygosity drops toward the heterozygosity of the disease locus ($m \rightarrow .1$). A typical example is at $r=5$ and $\delta=\frac{1}{2}\delta_{\max}$ where the table shows $P_t=.548$ at maximum marker heterozygosity ($m=.5$) and $P_t=.597$ or $.648$ for $m=.2$ or $.1$, respectively. The impact of such an increase in P_t can be understood by calculating χ^2_{tdt} for $P_t=.548$ ($m=.5$) and for $P_t=.597$ ($m=.2$) assuming a data set of 200 families each with two affected sibs. Based on the expression for $\frac{H}{F}$, I calculate the proportion of A/B parents to be .50 and .39 when $m=.5$ and $.2$, respectively. So for $m=.5$, there would be $.5 \times 400 \times 2 = 400$ informative transmissions to affected offspring with transmissions of allele A totaling $.548 \times 400 = 219$, thus implying $\chi^2_{\text{tdt}} = \frac{38^2}{400} = 3.61$, $p<0.1$. For $m=.2$, there would be $.39 \times 400 \times 2 = 312$ informative transmissions of which $.597 \times 312 = 186$ would be transmissions of allele A yielding $\chi^2_{\text{tdt}} = \frac{60^2}{312} = 11.54$, $p<0.005$.

This example is typical, and highlights perhaps the most important finding of this paper; namely the importance of using bi-allelic markers with heterozygosity similar to that of a bi-allelic disease locus. Indeed, since a majority of susceptibility loci may be bi-allelic, the

judicious use of bi-allelic markers of both high, medium, and low heterozygosity may be crucial in order to initially detect and replicate linkages to loci conferring modest disease risk.

Best Mode:

Method for locating disease causing polymorphism using biallelic linkage analysis

Objective :To test, by association-based linkage analysis (e.g., by TDT), whether a disease-causing polymorphism is located on a particular chromosome (e.g., human chromosome 4) or within a particular subregion of that chromosome.

PART 1 - Steps in conducting the association-based linkage test

Step 1

To conduct the test, first divide the chromosome or subregion of interest into segments that are short enough that polymorphisms within each segment are likely to be in linkage disequilibrium with each other. The division of a chromosome or subregion of interest into "segments" is conceptual (*not* physical) and is based on chromosomal maps such as those provided by the Whitehead Institute or Marshfield Foundation for Biomedical Research. Although disequilibrium has been observed in Finnish populations between polymorphisms that are 7 to 10 centimorgans (cM) apart, the chromosomal segments for searching for disease-causing polymorphisms in more genetically heterogeneous populations should be less than 1 cM long (e.g., 250,000 base pairs long). These chromosomal segments might or might not overlap each other (i.e., share some of their length in common); but the set of chromosomal segments should completely cover the entire chromosome or entire subregion of interest, so that a disease-causing polymorphism located anywhere on the chromosome or anywhere in the subregion of interest will be detected by the test.

Step 2

It is well known that increased disequilibrium between a marker and linked disease locus increases evidence for linkage provided by association-based linkage tests such as the TDT. However, what has not been recognized is that the specific allele frequencies of the marker locus can also have an enormous impact on the strength of evidence for linkage. I

showed this by analyzing equation 2 for P_t . Specifically, when a bi-allelic marker is in linkage disequilibrium with a bi-allelic disease locus, the strength of evidence for linkage provided by the TDT is *greatly* increased if the bi-allelic marker and bi-allelic disease locus have similar allele frequencies.

This phenomenon is illustrated by Table 2 and explained above. For example, suppose as noted above, that the susceptibility allele ("allele D") of a bi-allelic disease locus has a frequency of 0.1 and further suppose that the disease locus is in half-maximal positive disequilibrium with a bi-allelic marker ($\delta = \frac{1}{2} \delta_{\max}$). As noted above, χ^2_{TDT} will equal only 3.61 ($p < 0.1$) if the frequency of the less common marker allele is 0.5; but if the frequency of the less common marker allele is 0.2 (and hence much closer to the frequency of allele D) then χ^2_{TDT} will equal 11.54, thus providing much stronger evidence for linkage ($p < 0.005$).

Therefore, in searching for association-based linkage to a bi-allelic disease locus within each of the aforementioned chromosomal segments (see step 1), it is crucial to identify and test (e.g., by TDT) bi-allelic markers within each segment that have a broad range of allele frequencies. An unidentified bi-allelic disease locus could have allele frequencies close to 0.5/0.5, 0.4/0.6, 0.3/0.7, 0.2/0.8, 0.1/0.9 or below 0.1/above 0.9; hence, it is crucial to test bi-allelic markers with frequencies near 0.5/0.5 and near 0.1/0.9 as well as test others with allele frequencies that fall at regular increments between the extremes of 0.5/0.5 and 0.1/0.9. By testing bi-allelic markers with a broad range of allele frequencies that are spaced at regular intervals between 0.5/0.5 and 0.1/0.9, one is assured of testing some bi-allelic markers whose two allele frequencies are reasonably close to the allele frequencies of an unknown bi-allelic disease locus.

Thus, for step 2, within each chromosomal segment, subsets of bi-allelic markers should be identified. Each subset contains only bi-allelic markers having approximately the same allele frequencies. For example, subset A contains only markers whose less common allele has a population frequency of about 0.1. Similarly, subsets B, C, D, and E contain only bi-allelic markers whose less common allele has a frequency of approximately 0.2, 0.3, 0.4, and 0.5, respectively. In other versions of the invention the number of subsets is greater or less than five, and the approximate allele frequency of the less common bi-allele of subsets is other than about 0.1, 0.2, 0.3, 0.4 or 0.5 and is expected to be more than one decimal long since allele frequencies from real populations are rarely round numbers. However, the crucial point is that each subset should contain only bi-allelic markers belonging to one chromosomal segment and the frequency of the less common allele of each subset member should be

approximately the same (i.e., the *difference* between the frequencies of the less common allele of any two subset members should not exceed 0.15). Also crucial, as I emphasized above, is that the *group* of subsets for each chromosomal segment represent frequencies near the extremes of 0.5/0.5 and 0.1/0.9 as well as represent bi-allele frequencies between these two extremes that are approximately evenly spaced as *illustrated* by the group of subsets referred to above as A, B, C, D and E.

Step 3

In step 2, I described the importance of testing subsets of bi-allelic markers having approximately the same frequencies for their two alleles. Here I further delineate the characteristics of the markers that should be chosen for each subset by noting why it is important that each subset contain more than one bi-allelic marker. Even though a particular bi-allelic marker has allele frequencies that are similar to those of a closely linked bi-allelic disease locus, the marker may not be in strong positive disequilibrium with the disease locus. If disequilibrium is minimal, the marker will not show strong evidence for linkage under the TDT or any other association-based linkage test, *even though the bi-allelic marker and disease locus have similar allele frequencies*.

Hence, it is important that each subset contain multiple bi-allelic markers so that there is increased likelihood that at least one of the markers will be in reasonably strong disequilibrium with a closely linked bi-allelic disease locus. Beyond the cardinal criterion that all bi-allelic markers in a subset have similar allele frequencies, an additional criteria for selecting markers to belong to a subset is that the chosen bi-allelic markers *should not be in extreme positive disequilibrium with each other*.

The reason for this is as follows: According to standard usage, the disequilibrium coefficient (δ) is defined by the equation $\delta = f(AB) - f(A)f(B)$ where $f(A)$ and $f(B)$ may be defined as the frequencies of the less common allele (denoted A and B) of two bi-allelic loci belonging to the same subset and $f(AB)$ is the population frequency of the AB haplotype. Since the two markers belong to the same subset, we may assume that $f(A)=f(B)=q$; hence the maximum positive value of δ (denoted δ_{\max}) is $\delta=q-q^2$. This maximum positive δ value (i.e. maximum "positive disequilibrium") occurs if every chromosome that carries allele A also carries allele B, and if every chromosome that carries allele not-A also carries allele not-B. Hence, when two bi-allelic markers with similar allele frequencies are in extreme positive disequilibrium with each other (i.e., δ is approximately equal to δ_{\max}), the two loci provide

the nearly identical information with respect to their linkage and association with a third polymorphism such as a disease locus. Hence one of the two bi-allelic markers would provide no additional information and its inclusion in the subset would not increase the likelihood of detecting linkage and association to a nearby disease locus.

Therefore, bi-allelic markers belonging to the same chromosomal segment and subset should not only have similar allele frequencies, the δ value between *each pair* of bi-allelic markers in the same subset should be substantially less than $\delta_{\max} = q - q^2$. This assures that every bi-allelic polymorphism belonging to the subset provides much new (i.e. non-redundant) information about linkage and association to any nearby bi-allelic disease locus: thus testing each bi-allelic marker in the subset would increase the likelihood of detecting linkage to a disease locus.

Step4: Test for linkage

To test for (association-based) linkage to a bi-allelic disease locus, each bi-allelic marker in each subset from each chromosomal segment is tested *individually* by using the TDT, AFBAC method or other family-based linkage test. To conduct these tests for a particular marker, members of nuclear families (most especially parents, and any children who manifest disease) are genotyped at the marker being tested and the genotypes are then evaluated according to the TDT, AFBAC method or other family-based linkage/association test (for description of TDT and AFBAC, see Spielman et al, Am J of Human Genetics 52:506-516 (1993) and Thomson, Am J Human Genetics 57:487-498 (1995)). Alternatively, linkage and association is tested for each marker in each subset from each segment by genotyping individuals with disease and related or unrelated normal controls at each marker to be tested. (End of best mode example)

Further Information

(Step 3 is not essential for the operation or utility of this version of the invention. In this best mode example, the least common allele frequency subrange 0.1 to 0.5 is used. In versions of the invention similar to the best mode, versions of the invention are operable and have utility for any subrange of the least common allele frequency range 0 to 0.5. In addition, rather than genotyping DNA from single individuals in step 4, in some versions of the invention each marker in each subset from each segment is tested for association with disease by evaluating DNA from pooled samples.)

PART 2 - Physical implementation of the above test

Silicon chips or glass slides with arrays of oligonucleotides for testing bi-allelic markers

Companies like Affymetrix[™] are using silicon chips or glass slides to genotype DNA from one individual at thousands of bi-allelic markers. Each silicon chip or glass slide is divided into a grid or 2-dimensional matrix that contains thousands of cells. To the surface of each cell is attached multiple copies of a unique oligonucleotide whose sequence is complementary (type (1)) to one of the two alleles of a particular bi-allelic marker. Thus, DNA from an individual who carries the allele hybridizes to the cell with substantially greater affinity than does the alternate bi-allele. The degree of hybridization generates a signal and enables the genotype of the individual to be inferred for that particular bi-allelic polymorphism [i.e., the individual is homozygous (++) , heterozygous (+-), or homozygous (--)]. In some applications, it is crucial to attach oligonucleotides corresponding to each allele of a bi-allelic polymorphism in adjacent cells so that the relative (i.e. local) intensity of hybridization in the adjacent cells can be compared, thus facilitating inference of the individual's correct genotype (++, +-, or --).

In using this silicon chip or glass slide technology to test for linkage and association, the ideas detailed in PART 1 indicate how the oligonucleotides that are attached to the cells of the silicon chip or glass slide should be chosen. To scan a particular chromosome or chromosomal region for a bi-allelic disease locus, the chromosome or chromosomal region should be subdivided into segments as described in Step 1 above. For each segment, subsets of bi-allelic markers having the properties detailed in PART 1 above should be identified. The DNA of select individuals (see "Test for linkage" - above) should then be assayed at each bi-allelic marker in every chromosomal segment and in every subset of markers belonging to the segment. This would be accomplished by attaching an oligonucleotide corresponding to one of the marker's two alleles to a particular (i.e. known) cell on the silicon chip or slide. To enhance assignment of accurate genotypes, it may also be advisable to attach an oligonucleotide corresponding to the second allele of the bi-allelic marker in an adjacent cell as mentioned in the previous paragraph.

Industrial Applicability

Versions of the present invention are useful for locating trait causing genes and polymorphisms such as human disease genes and polymorphisms. Versions of the invention could be used to find the cure for human disease. The making and use of versions of the invention should be clear to a person of skill in the art after reading the description.

Scope of the Invention

While the description contains many specificities, these should not be construed as limitations on the scope of the invention, but rather as exemplifications of versions of the invention.

Accordingly the scope of the invention should be determined not by the specific versions described alone, but also by the claims and their legal equivalents and also by any future claims drawn to the invention and future descriptions of versions of the invention.

Notes:

The reader's attention is directed to the following papers which are open to the public and are herein incorporated by reference: (1) McGinnis, Ewens & Spielman, Genetic Epidemiology 1995 ; 12(6) : 637-40. (2) RE McGinnis Annals of Human Genetics vol 62, pp. 159-179, 1998. The papers in the endnotes below are incorporated herein by reference.

ⁱ Weighing DNA for Fast Genetic Diagnosis, Science, March 27, 1998, vol. 279, pp. 2044-2045.

ⁱⁱ Spielman, R.S. and Ewens, W.J. The TDT and Other Family-Based Tests for Linkage Disequilibrium and Association, American Journal of Human Genetics, 59: 983-989, 1996.

ⁱⁱⁱ "Mathematical Theory of Optimization" The New Encyclopedia Britannica, 15th edition, vol. 25, pp. 217-221.

^{iv} American Journal of Human Genetics, vol. 57: 439-454, 1995.

^v American Journal of Human Genetics, vol. 58: 1347-1363, 1996.

^{vi} Human Heredity, vol. 44, pp. 225-237, 1994.

^{vii} Human Heredity, vol. 46, pp. 226-235, 1996.

^{viii} Accessing Genetic Information with High-Density DNA Arrays, Mark Chee, et al. Science, vol 274, Oct. 25, 1996, pp. 610 - 614.

^{ix} Large Scale Identification, Mapping, and Genotyping of Single-Nucleotide Polymorphisms in the Human Genome, Wang, et al., Science, May 15, 1998, vol 280, pp. 1077-1081.

^x (1) Schuster, H. et al (1995) Nature Genetics, 13(1) : 98 - 100.

(2) Gyapay, G. et al (1994) Nature Genetics, 7: 246-339.

^{xi} Some versions of oligonucleotide technology and its uses to obtain genetic information are included in the following papers:

(1) Accessing Genetic Information with High-Density DNA Arrays, Mark Chee, et al. Science, vol 274, Oct. 25, 1996, pp. 610 - 614.

(2) Genetic analysis of amplified DNA with immobilized sequence-specific oligonucleotide probes, Saiki, et al. Proc Natl Acad Sci USA vol 86, pp. 6230-6234.

(3) Allele-specific enzymatic amplification of β -globin genomic DNA for diagnosis of sickle cell anemia, Wu, et al., Proc Natl Acad Sci USA vol 86 pp 2757-2760.

(4) Automated DNA diagnostics using an Elisa-based oligonucleotide ligation assay, Nickerson, et al., Proc Natl Acad Sci USA vol 87, pp. 8923-8927.

(5) Genetic analysis of amplified DNA with immobilized sequence specific oligonucleotide probes, Saiki, et al., Proc Natl Acad Sci USA vol 86 pp 6230 - 6234.

^{xii} Padlock Probes: Circularizing Oligonucleotides for Localized DNA Detection, Science, Sept. 30, 1994, vol. 265, pp. 2085-2088.

^{xiii} SNP attack on complex traits, Nature Genetics, Nov. 1998, vol. 20 no. 3, pp. 217-218.

Claims

What is claimed:

1. An invention as described in the description.
2. A process for identifying one or more bi-allelic markers linked to a bi-allelic genetic characteristic gene in a species of creatures, comprising the steps of :
 - a) choosing two or more bi-allelic covering markers so that a CL-F region is systematically covered by the two or more covering markers;
 - b) choosing a statistical linkage test based on allelic association for each covering marker;
 - c) choosing a sample of individuals for each covering marker ;
 - d) obtaining genotype data/sample allele frequency data for each covering marker and the sample chosen for each covering marker, and obtaining phenotype status data for the genetic characteristic for each individual in the sample chosen for each covering marker;
 - e) calculating evidence for linkage between each covering marker and the gene using the statistical linkage test based on allelic association chosen for each covering marker and the genotype data/sample allele frequency data for each covering marker and using the phenotype status data for the genetic characteristic for each individual in the sample chosen for each covering marker obtained in d); and
 - f) identifying those covering markers as linked to the genetic characteristic gene which show evidence for linkage based on the calculations of step e.
3. A process as in claim 2, wherein the CL-F region is N covered to within a CL-F distance δ by the two or more bi-allelic covering markers, wherein δ is equal to about $[\delta_{CL}, 0.25]$ or the equivalent thereof, δ_{CL} is equal to the largest chromosomal length, computed by any method, for which linkage disequilibrium has been observed between any polymorphisms in any population of the species, N is an integer greater than or equal to 1.
4. A process as in claim 2, wherein the CL-F region is N covered to within a CL-F distance δ by the two or more bi-allelic covering markers, wherein δ is less than or equal to about $[\delta_{CL}, \delta_F]$ or the equivalent thereof, δ_{CL} is equal to the largest chromosomal length, computed by any method, for which linkage disequilibrium has been observed between any polymorphisms in any population of the species, δ_F is equal to 0.25, N is an integer greater than or equal to 1, the CL-F region comprising a CL-F matrix and there being N or more covering markers in each cell of the matrix, each cell of the matrix being of length

- 1 L_{MC} and width W_{MC} , and there being three or more columns in the matrix and one or more rows in the
 2 matrix and L_{MC} being less than or equal to δ_{CL} , and W_{MC} being less than or equal to 0.25.
- 3 5. A process as in claim 4, wherein W_{MC} is less than 0.15, wherein R_M is the number of rows in the
 4 matrix, wherein C_M is the number of columns in the matrix, wherein A_M is the average chromosomal
 5 intermarker distance of the covering markers in the matrix, wherein M is the total number of covering
 6 markers in the matrix, wherein (1) M is less than or equal to $N(R_M + 1)C_M$ when less than 10% of the
 7 chromosomal intermarker distances of the covering markers in the matrix are greater than $(A_M/2)$ and
 8 (2) wherein M has no upper bound when greater than or equal to 10% of the chromosomal intermarker
 9 distances of the covering markers in the matrix are greater than $(A_M/2)$ and wherein at least one
 10 covering marker in the matrix has a least common allele frequency less than 0.4, wherein the species is
 11 human being, and wherein each covering marker is an SNP.
- 12 6. A process for localizing a bi-allelic genetic characteristic gene in a species of creatures to a CL-F
 13 location, comprising the steps of:
 14
- 15 a)choosing two or more bi-allelic covering markers so that a CL-F region is systematically covered by
 16 the two or more covering markers, the CL-F region comprising a CL-F matrix and there being N or more
 17 covering markers in each cell of the matrix, each cell of the matrix being of length L_{MC} and width W_{MC} ,
 18 and there being three or more columns in the matrix and one or more rows in the matrix and L_{MC} being
 19 less than or equal to the largest chromosomal length, computed by any method, for which linkage
 20 disequilibrium has been observed between any polymorphisms in any population of the species, and
 21 W_{MC} being less than or equal to 0.25, N being an integer greater than or equal to 1;
 22
- 23 b)choosing a statistical linkage test based on allelic association for each covering marker;
 24
- 25 c)choosing a sample of individuals for each covering marker ;
 26
- 27 d)obtaining genotype data/sample allele frequency data for each covering marker and the sample
 28 chosen for each covering marker, and obtaining phenotype status data for the genetic characteristic for
 29 each individual in the sample chosen for each covering marker;
 30
- 31 e)calculating evidence for linkage between each covering marker and the gene using the statistical
 32 linkage test based on allelic association chosen for each covering marker and the genotype
 33 data/sample allele frequency data for each covering marker and using the phenotype status data for the
 34 genetic characteristic for each individual in the sample chosen for each covering marker obtained in d);
 35 and
 36
- 37 f)localizing the gene to the CL-F location of one or more markers that show evidence for linkage based
 38 on the calculations of step e).
- 39 7.A process for obtaining genotype data/sample allele frequency data for each bi-allelic marker of a
 40 group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of

1 a sample, (each individual being a member of one species), comprising:

2
3 a) determining information on the presence or absence of each allele of each bi-allelic marker of a
4 group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of
5 a sample, a CL-F region being systematically covered by the two or more bi-allelic covering markers,
6 the CL-F region comprising a CL-F matrix and there being N or more covering markers in each cell of
7 the matrix, each cell of the matrix being of length L_{MC} and width W_{MC} , and there being three or more
8 columns in the matrix and one or more rows in the matrix and L_{MC} being less than or equal to the largest
9 chromosomal length, computed by any method, for which linkage disequilibrium has been observed
10 between any polymorphisms in any population of the species, and W_{MC} being less than or equal to
11 0.25, N being an integer greater than or equal to 1; and

12
13 b) transforming the information of step a) into genotype data/sample allele frequency data for each
14 marker of the group.

15 8. A process as in any one of claims 2, 3, 4, 5, 6, or 7, wherein the least common allele frequency of
16 one or more of the covering markers is less than 0.4, wherein the species is human being, and wherein
17 each covering marker is an SNP.

18 9. A process as in any one of claims 2, 3, 4, 5, 6, or 7, wherein the species is human being, and
19 wherein each covering marker is an SNP, and wherein (1) the least common allele frequency of one or
20 more of the covering markers is less than 0.3 and the chromosomal intermarker distances of the
21 covering markers are approximately equal and the average chromosomal intermarker distance of the
22 covering markers is greater than 10 cM or an equivalent thereof, or wherein (2) the least common allele
23 frequency of one or more of the covering markers is less than 0.2 and the chromosomal intermarker
24 distances of the covering markers are approximately equal and the average chromosomal intermarker
25 distance of the covering markers is greater than 1 cM or an equivalent thereof.

26 10. An apparatus for identifying bi-allelic markers linked to a bi-allelic genetic characteristic gene in a
27 species of creatures, comprising :

28
29 a) means for choosing two or more bi-allelic covering markers so that a CL-F region is systematically
30 covered by the two or more covering markers, the CL-F region comprising a CL-F matrix and there
31 being N or more covering markers in each cell of the matrix, each cell of the matrix being of length L_{MC}
32 and width W_{MC} , and there being three or more columns in the matrix and one or more rows in the matrix
33 and L_{MC} being less than or equal to the largest chromosomal length, computed by any method, for
34 which linkage disequilibrium has been observed between any polymorphisms in any population of the
35 species, and W_{MC} being less than or equal to 0.25, N being an integer greater than or equal to 1;

36
37 b) means for choosing a statistical linkage test based on allelic association for each covering marker;

38
39 c) means for choosing a sample of individuals for each covering marker ;

d) means for obtaining genotype data/sample allele frequency data for each covering marker and the sample chosen for each covering marker, and for obtaining phenotype status data for the genetic characteristic for each individual in the sample chosen for each covering marker;

e) means for calculating evidence for linkage between each covering marker and the gene using the statistical linkage test based on allelic association chosen for each covering marker and the genotype data/sample allele frequency data for each covering marker and using the phenotype status data for the genetic characteristic for each individual in the sample chosen for each covering marker obtained in d); and

f) means for identifying those covering markers as linked to the gene which show evidence for linkage based on the calculations by means e).

11. An apparatus for obtaining genotype data/sample allele frequency data for each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of a sample, (each individual being a member of one species), comprising:

a) means for determining information on the presence or absence of each allele of each bi-allelic marker of a group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of the sample, a CL-F region being systematically covered by the two or more bi-allelic covering markers, the CL-F region comprising a CL-F matrix and there being N or more covering markers in each cell of the matrix. each cell of the matrix being of length L_{MC} and width W_{MC} , and there being three or more columns in the matrix and one or more rows in the matrix and L_{MC} being less than or equal to the largest chromosomal length, computed by any method, for which linkage disequilibrium has been observed between any polymorphisms in any population of the species, and W_{MC} being less than or equal to 0.25, N being an integer greater than or equal to 1; and

b) means for transforming the information of step a) into genotype data/sample allele frequency data for each marker of the group.

12. An apparatus as in claim 10 or claim 11, wherein the least common allele frequency of one or more of the covering markers is less than 0.4, wherein the species is human being, and wherein each covering marker is an SNP.

13. An apparatus as in claim 10 or claim 11, wherein the species is human being, and wherein each covering marker is an SNP, and wherein (1) the least common allele frequency of one or more of the covering markers is less than 0.3 and the chromosomal intermarker distances of the covering markers are approximately equal and the average chromosomal intermarker distance of the covering markers is greater than 10 cM or an equivalent thereof, or wherein (2) the least common allele frequency of one or

more of the covering markers is less than 0.2 and the chromosomal intermarker distances of the covering markers are approximately equal and the average chromosomal intermarker distance of the covering markers is greater than 1 cM or an equivalent thereof .

14. The use of one or more copies of a set of oligonucleotides to determine genotype data/sample allele frequency data for each bi-allelic marker of a group of two or more bi-allelic covering markers for one or more individuals, (each individual being a member of one species), wherein the group of covering markers systematically cover a CL-F region, the CL-F region comprising a CL-F matrix and there being N or more covering markers in each cell of the matrix, each cell of the matrix being of length L_{MC} and width W_{MC} , and there being three or more columns in the matrix and one or more rows in the matrix and L_{MC} being less than or equal to the largest chromosomal length, computed by any method, for which linkage disequilibrium has been observed between any polymorphisms in any population of the species, and W_{MC} being less than or equal to 0.25, N being an integer greater than or equal to 1.

15. The use in claim 14, wherein the least common allele frequency of one or more of the covering markers is less than 0.4, wherein the species is human being, and wherein each covering marker is an SNP.

16. The use in claim 14, wherein the species is human being, and wherein each covering marker is an SNP, and wherein (1) the least common allele frequency of one or more of the covering markers is less than 0.3 and the chromosomal intermarker distances of the covering markers are approximately equal and the average chromosomal intermarker distance of the covering markers is greater than 10 cM or an equivalent thereof, or wherein (2) the least common allele frequency of one or more of the covering markers is less than 0.2 and the chromosomal intermarker distances of the covering markers are approximately equal and the average chromosomal intermarker distance of the covering markers is greater than 1 cM or an equivalent thereof .

17. One or more copies of a set of oligonucleotides, the set of oligonucleotides being complementary to a group of two or more bi-allelic covering markers (of a species), wherein the group of covering markers systematically cover a CL-F region, the CL-F region comprising a CL-F matrix and there being N or more covering markers in each cell of the matrix, each cell of the matrix being of length L_{MC} and width W_{MC} , and there being three or more columns in the matrix and one or more rows in the matrix and L_{MC} being less than or equal to the largest chromosomal length, computed by any method, for which linkage disequilibrium has been observed between any polymorphisms in any population of the species, and W_{MC} being less than or equal to 0.25, N being an integer greater than or equal to 1.

18. One or more copies of a set of oligonucleotides as in claim 17, wherein the least common allele frequency of one or more of the covering markers is less than 0.4, wherein the species is human being, and wherein each covering marker is an SNP.

19. One or more copies of a set of oligonucleotides as in claim 17, wherein the species is human being, and wherein each covering marker is an SNP, and wherein (1) the least common allele frequency of one or more of the covering markers is less than 0.3 and the chromosomal intermarker distances of the covering markers are approximately equal and the average chromosomal intermarker distance of the covering markers is greater than 10 cM or an equivalent thereof, or wherein (2) the least common allele

- 1 frequency of one or more of the covering markers is less than 0.2 and the chromosomal intermarker
2 distances of the covering markers are approximately equal and the average chromosomal intermarker
3 distance of the covering markers is greater than 1 cM or an equivalent thereof .
- 4 20. One or more copies of a set of oligonucleotides as in claim 17, wherein less than or equal to R pairs
5 of the covering markers are redundant within CL-F distance D, wherein D is equal to $[L_{MC}, W_{MC}]$,
6 wherein R is an integer greater than or equal to 0 and less than or equal to about $N(N-1)/2$.
- 7 21. An apparatus as in any one of claims 10, 11, 12, or 13, wherein the apparatus comprises a
8 computer, the computer being supplied with proper data and instructions.
- 9 22. A process for localizing a bi-allelic genetic characteristic gene in a species of creatures to a CL-F
10 location, comprising the steps a), b), c), d) and e) of the process in any one of claims 2, 3, 4, 5, and
11 further comprising the step of:
- 12 f)localizing the gene to the CL-F location of one or more markers that show evidence for linkage based
13 on the calculations of step e).
- 14 23. A process as in claim 22 wherein the least common allele frequency of one or more of the covering
15 markers is less than 0.4. wherein the species is human being, and wherein each covering marker is an
16 SNP.
- 17 24 A process as in 22 wherein the species is human being, and wherein each covering marker is an
18 SNP, and wherein (1) the least common allele frequency of one or more of the covering markers is less
19 than 0.3 and the chromosomal intermarker distances of the covering markers are approximately equal
20 and the average chromosomal intermarker distance of the covering markers is greater than 10 cM or an
21 equivalent thereof, or wherein (2) the least common allele frequency of one or more of the covering
22 markers is less than 0.2 and the chromosomal intermarker distances of the covering markers are
23 approximately equal and the average chromosomal intermarker distance of the covering markers is
24 greater than 1 cM or an equivalent thereof .
- 25 25. A process as in any one of claims 2, 3, 4, 5, 6, 21, 22, 23, or 24 wherein the process comprises a
26 computer program.
- 27

AMENDED CLAIMS

[received by the International Bureau on 26 July 1999 (26.07.99);
original claim 1 amended; original claims 2- 25 cancelled;
new claims 2-88 added; (19 pages)]

4 What is claimed:

5 1. An invention substantially as described in the description.

6 2. An invention substantially as described and illustrated in the description.

7 3. A process for identifying one or more bi-allelic markers linked to a bi-allelic genetic characteristic
8 gene in a species of creatures, comprising the steps of :

9

10 a)choosing two or more bi-allelic covering markers so that a CL-F region is systematically covered by
11 the two or more covering markers, the CL-F region being a collection of points on a two-dimensional
12 plane, the two-dimensional plane having the two orthogonal dimensions of chromosomal location and
13 least common allele frequency;

14

15 b)choosing a statistical linkage test based on allelic association for each covering marker;

16

17 c)choosing a sample of individuals for each covering marker ;

18

19 d)obtaining genotype data/sample allele frequency data for each covering marker and the sample
20 chosen for each covering marker, and obtaining phenotype status data for the genetic characteristic for
21 each individual in the sample chosen for each covering marker;

22

23 e)calculating evidence for linkage between each covering marker and the gene using the statistical
24 linkage test based on allelic association chosen for each covering marker and the genotype
25 data/sample allele frequency data for each covering marker and using the phenotype status data for the
26 genetic characteristic for each individual in the sample chosen for each covering marker obtained in d);
27 and

28

29 f)identifying those covering markers as linked to the genetic characteristic gene which show evidence
30 for linkage based on the calculations of step e.

31 4. A process as in claim 3, wherein the CL-F region is N covered to within a CL-F distance δ by the two
32 or more bi-allelic covering markers, so that each point in the region is within the CL-F distance δ of N or
33 more of the covering markers, wherein δ is equal to about $[\delta_{CL}, 0.25]$ or the equivalent thereof, δ_{CL} is
34 equal to the largest chromosomal length, computed by any method, for which linkage disequilibrium has
35 been observed between any polymorphisms in any population of the species, N is an integer greater
36 than or equal to 1.

37 5. A process as in claim 4, wherein the CL-F region comprises a CL-F matrix, the sum of the number of
38 columns and rows in the matrix being greater than or equal to three, each cell of the matrix being of
39 length L_{MC} and width W_{MC} , and L_{MC} being less than or equal to about δ_{CL} , and W_{MC} being less than or
40 equal to about 0.25, δ_{CL} is equal to the largest chromosomal length, computed by any method, for which

1 linkage disequilibrium has been observed between any polymorphisms in any population of the species,
2 there being N or more covering markers in each cell of the matrix and N is an integer greater than or
3 equal to 1.

4 6. A process as in claim 5, wherein the covering markers are substantially nonevenly distributed across
5 a chromosome or a chromosomal segment.

6 7. A process as in claim 5, wherein the covering markers are substantially evenly distributed across a
7 chromosome or a chromosomal segment, and wherein the least common allele frequency of one or
8 more markers is less than 0.4.

9 8. A process as in claim 5, wherein the covering markers are substantially evenly distributed across a
10 chromosome or a chromosomal segment; and wherein there is a subgroup of one or more of the
11 covering markers, and each of the markers in the subgroup is chosen without substantial preference for
12 the least common allele frequency of each of the markers in the subgroup being close to 0.5.

13 9. A process as in claim 5, wherein the covering markers are substantially evenly distributed across a
14 chromosome or a chromosomal segment, wherein (1) the average chromosomal intermarker distance
15 of the covering markers is greater than 2 cM or the equivalent thereof and the least common allele
16 frequency of one or more of the covering markers is less than 0.3 or wherein (2) the least common
17 allele frequency of one or more of the covering markers is less than 0.2.

18 10. A process as in claim 5, wherein the covering markers are substantially evenly distributed across a
19 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
20 the covering markers is less than or equal to 2 cM or the equivalent thereof, and wherein the conditional
21 probability the covering markers were chosen essentially randomly from substantially the known set of
22 bi-allelic markers with least common allele frequencies between 0.2 inclusive and 0.5 inclusive is less
23 than about 10 percent; wherein the conditional probability is substantially conditional on (1) the
24 approximate chromosomal distribution of the covering markers, (2) the marker type of each covering
25 marker and (3) there being N or more covering markers in each cell of the matrix.

26 11. A process as in claim 5, wherein the covering markers are substantially evenly distributed across a
27 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
28 the covering markers is greater than 2 cM or the equivalent thereof; and wherein the conditional
29 probability the covering markers were chosen essentially randomly from substantially the known set of
30 bi-allelic markers with least common allele frequencies between 0.3 inclusive and 0.5 inclusive is less
31 than about 10 percent; wherein the conditional probability is substantially conditional on (1) the
32 approximate chromosomal distribution of the covering markers, (2) the marker type of each covering
33 marker and (3) there being N or more covering markers in each cell of the matrix.

34 12. A process as in claim 10, wherein the chromosome or the chromosomal segment consists
35 essentially of a set of nonoverlapping chromosome segments of substantially equal length, and wherein
36 one and only one covering marker is located on each of 80 percent or more of the chromosome
37 segments of the set, and wherein zero or two and only two covering markers are located on each of 20
38 percent or less of the chromosome segments of the set, and wherein each chromosome segment with
39 zero covering markers located thereon is bordered only by chromosome segments with one or two
40 covering markers located thereon, and wherein each chromosome segment with two covering markers

1 located thereon is bordered only by chromosome segments with one or zero covering markers located
2 thereon; and wherein the conditional probability the covering markers were chosen essentially randomly
3 from substantially the known set of bi-allelic markers with least common allele frequencies between 0.2
4 inclusive and 0.5 inclusive is less than about 10 percent; wherein the conditional probability is
5 substantially conditional on (1) the chromosomal distribution of the covering markers on the
6 chromosome segments of the set, (2) the marker type of each covering marker and (3) there being N or
7 more covering markers in each cell of the matrix.

8 13. A process as in claim 11, wherein the chromosome or the chromosomal segment consists
9 essentially of a set of nonoverlapping chromosome segments of substantially equal length, and wherein
10 one and only one covering marker is located on each of 80 percent or more of the chromosome
11 segments of the set, and wherein zero or two and only two covering markers are located on each of 20
12 percent or less of the chromosome segments of the set, and wherein each chromosome segment with
13 zero covering markers located thereon is bordered only by chromosome segments with one or two
14 covering markers located thereon, and wherein each chromosome segment with two covering markers
15 located thereon is bordered only by chromosome segments with one or zero covering markers located
16 thereon; and wherein the conditional probability the covering markers were chosen essentially randomly
17 from substantially the known set of bi-allelic markers with least common allele frequencies between 0.3
18 inclusive and 0.5 inclusive is less than about 10 percent; wherein the conditional probability is
19 substantially conditional on (1) the chromosomal distribution of the covering markers on the
20 chromosome segments of the set, (2) the marker type of each covering marker and (3) there being N or
21 more covering markers in each cell of the matrix.

22 14. A process as in claim 5, wherein the covering markers are substantially evenly distributed across a
23 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
24 the covering markers is less than or equal to 2 cM or the equivalent thereof, and wherein collection C is
25 essentially the collection of known groups of bi-allelic markers with least common allele frequencies
26 between 0.2 inclusive and 0.5 inclusive that are substantially similar to the covering markers as a
27 group; wherein a group of bi-allelic markers is a member of collection C if and only if the group
28 substantially meets criteria (1), (2), (3) and (4): (1) each marker in the group is chosen from
29 substantially the known set of bi-allelic markers with least common allele frequencies between 0.2
30 inclusive and 0.5 inclusive, (2) the number of markers in the group is the same as the number of
31 covering markers, (3) the chromosomal distribution of the group of markers and the covering markers is
32 substantially similar, and (4) the marker type of each group marker and the covering marker with
33 substantially the same chromosomal location is the same; wherein a group that is a member of
34 collection C substantially meets criterion (5) if and only if (5) there are N or more of the group markers
35 in each cell of the matrix; wherein P is essentially the proportion of groups in collection C that meet
36 criterion (5); wherein P is less than about 90 percent.

37 15. A process as in claim 5, wherein the covering markers are substantially evenly distributed across a
38 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
39 the covering markers is greater than 2 cM or the equivalent thereof, and wherein collection C is
40 essentially the collection of known groups of bi-allelic markers with least common allele frequencies

1 between 0.3 inclusive and 0.5 inclusive that are substantially similar to the covering markers as a
2 group; wherein a group of bi-allelic markers is a member of collection C if and only if the group
3 substantially meets criteria (1), (2), (3) and (4): (1) each marker in the group is chosen from
4 substantially the known set of bi-allelic markers with least common allele frequencies between 0.3
5 inclusive and 0.5 inclusive, (2) the number of markers in the group is the same as the number of
6 covering markers, (3) the chromosomal distribution of the group of markers and the covering markers is
7 substantially similar, and (4) the marker type of each group marker and the covering marker with
8 substantially the same chromosomal location is the same; wherein a group that is a member of
9 collection C substantially meets criterion (5) if and only if (5) there are N or more of the group markers
10 in each cell of the matrix; wherein P is essentially the proportion of groups in collection C that meet
11 criterion (5); wherein P is less than about 90 percent.

12 16. A process as in claim 5, wherein the covering markers are substantially evenly distributed across a
13 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
14 the covering markers is less than or equal to 2 cM or the equivalent thereof; wherein the chromosome
15 or the chromosomal segment consists essentially of a set of nonoverlapping chromosome segments of
16 substantially equal length, and wherein one and only one covering marker is located on each of 80
17 percent or more of the chromosome segments of the set, and wherein zero or two and only two
18 covering markers are located on each of 20 percent or less of the chromosome segments of the set,
19 and wherein each chromosome segment with zero covering markers located thereon is bordered only
20 by chromosome segments with one or two covering markers located thereon, and wherein each
21 chromosome segment with two covering markers located thereon is bordered only by chromosome
22 segments with one or zero covering markers located thereon; wherein collection D is essentially the
23 collection of known groups of bi-allelic markers with least common allele frequencies between 0.2
24 inclusive and 0.5 inclusive that are substantially similar to the covering markers as a group; wherein a
25 group of bi-allelic markers is a member of collection D if and only if the group substantially meets
26 criteria (1), (2), and (3): (1) each marker in the group is chosen from substantially the known set of bi-
27 allelic markers with least common allele frequencies between 0.2 inclusive and 0.5 inclusive, (2) the
28 number of covering markers and the number of group markers located on each chromosome segment
29 of the set is the same, and (3) there is a group marker of the same type as each covering marker
30 located on the same chromosome segment of the set as each covering marker; wherein a group that is
31 a member of collection D substantially meets criterion (5) if and only if (5) there are N or more of the
32 group markers in each cell of the matrix; wherein P is essentially the proportion of groups in collection D
33 that meet criterion (5); wherein P is less than about 90 percent.

34 17. A process as in claim 5, wherein the covering markers are substantially evenly distributed across a
35 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
36 the covering markers is greater than 2 cM or the equivalent thereof; wherein the chromosome or the
37 chromosomal segment consists essentially of a set of nonoverlapping chromosome segments of
38 substantially equal length, and wherein one and only one covering marker is located on each of 80
39 percent or more of the chromosome segments of the set, and wherein zero or two and only two
40 covering markers are located on each of 20 percent or less of the chromosome segments of the set,

- 1 and wherein each chromosome segment with zero covering markers located thereon is bordered only
2 by chromosome segments with one or two covering markers located thereon, and wherein each
3 chromosome segment with two covering markers located thereon is bordered only by chromosome
4 segments with one or zero covering markers located thereon; wherein collection D is essentially the
5 collection of known groups of bi-allelic markers with least common allele frequencies between 0.3
6 inclusive and 0.5 inclusive that are substantially similar to the covering markers as a group; wherein a
7 group of bi-allelic markers is a member of collection D if and only if the group substantially meets
8 criteria (1), (2), and (3): (1) each marker in the group is chosen from substantially the known set of bi-
9 allelic markers with least common allele frequencies between 0.3 inclusive and 0.5 inclusive, (2) the
10 number of covering markers and the number of group markers located on each chromosome segment
11 of the set is the same, and (3) there is a group marker of the same type as each covering marker
12 located on the same chromosome segment of the set as each covering marker; wherein a group that is
13 a member of collection D substantially meets criterion (5) if and only if (5) there are N or more of the
14 group markers in each cell of the matrix; wherein P is essentially the proportion of groups in collection D
15 that meet criterion (5); wherein P is less than about 90 percent.
- 16 18. A process as in claim 4, wherein δ is less than or equal to about [1 cM, 0.15] or the equivalent
17 thereof.
- 18 19. A process as in claim 4, wherein (1) the covering markers are substantially nonevenly distributed
19 across a chromosome or a chromosomal segment or (2) wherein the covering markers are substantially
20 evenly distributed across a chromosome or a chromosomal segment, and wherein the least common
21 allele frequency of one or more markers is less than 0.4 or (3) wherein the covering markers are
22 substantially evenly distributed across a chromosome or a chromosomal segment; and wherein there is
23 a subgroup of one or more of the covering markers, and each of the markers in the subgroup is chosen
24 without substantial preference for the least common allele frequency of each of the markers in the
25 subgroup being close to 0.5.
- 26 20. A process as in any one of claims 3-19, wherein there is a group of covering markers, and the
27 markers in the group are a majority of the covering markers, and each marker in the group is an SNP,
28 or a bi-allelic marker equivalent formed only from one or more SNPs.
- 29 21. A process as in claim 5 wherein L_{MC} is less than or equal to about 250,000 bp or the equivalent
30 thereof, W_{MC} is less than or equal to about 0.15, wherein the species is human being, wherein the same
31 statistical linkage test based on allelic association is chosen for each covering marker in step b) and
32 wherein there is a group of covering markers, and the markers in the group are a majority of the
33 covering markers, and each marker in the group is an SNP, or a bi-allelic marker equivalent formed
34 only from one or more SNPs.
- 35 22. An apparatus for identifying bi-allelic markers linked to a bi-allelic genetic characteristic gene in a
36 species of creatures, comprising: means to practice each of the steps of a process as in any one of the
37 claims 3-21.
- 38 23. A process as in any one of claims 3-21, wherein the process comprises a computer program.
- 39 24. An apparatus as in claim 22, wherein the apparatus comprises a computer. the computer being
40 supplied with proper data and instructions.

1 25. A process for localizing a bi-allelic genetic characteristic gene in a species of creatures to a CL-F
2 location, comprising the steps of: any one of the processes in claims 3-21; further comprising: the step
3 f) localizing the gene to the CL-F location of one or more markers that show evidence for linkage based
4 on the calculations of step e).

5 26. A process as in claim 25, wherein the process comprises a computer program.

6 27. An apparatus for localizing a bi-allelic genetic characteristic gene in a species of creatures to a CL-F
7 location, comprising: means to practice each of the steps of a process as in claim 25.

8 28. An apparatus for localizing a bi-allelic genetic characteristic gene in a species of creatures to a CL-F
9 location as in claim 27, wherein the apparatus comprises a computer, the computer being supplied with
10 proper data and instructions.

11
12 29. A process for obtaining genotype data/sample allele frequency data for each bi-allelic marker of a
13 group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of
14 a sample, each individual in the sample being a member of the same species, comprising:

15
16 a) determining information on the presence or absence of each allele of each bi-allelic marker of a group
17 of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals of a
18 sample, a CL-F region being systematically covered by the two or more bi-allelic covering markers, the
19 CL-F region being a collection of points on a two-dimensional plane, the two-dimensional plane having
20 the two orthogonal dimensions of chromosomal location and least common allele frequency; and

21
22 b) transforming the information of step a) into genotype data/sample allele frequency data for each
23 marker of the group.

24
25 30. A process for obtaining genotype data/sample allele frequency data as in claim 29, wherein the CL-
26 F region is N covered to within a CL-F distance δ by the two or more bi-allelic covering markers, so that
27 each point in the region is within the CL-F distance δ of N or more of the covering markers, wherein δ is
28 equal to about $[\delta_{CL}, 0.25]$ or the equivalent thereof, δ_{CL} is equal to the largest chromosomal length,
29 computed by any method, for which linkage disequilibrium has been observed between any
30 polymorphisms in any population of the species, N is an integer greater than or equal to 1.

31 31. A process for obtaining genotype data/sample allele frequency data as in claim 30, wherein the CL-
32 F region comprises a CL-F matrix, the sum of the number of columns and rows in the matrix being
33 greater than or equal to three, each cell of the matrix being of length L_{MC} and width W_{MC} , and L_{MC} being
34 less than or equal to about δ_{CL} , and W_{MC} being less than or equal to about 0.25, δ_{CL} is equal to the
35 largest chromosomal length, computed by any method, for which linkage disequilibrium has been
36 observed between any polymorphisms in any population of the species, there being N or more covering
37 markers in each cell of the matrix and N is an integer greater than or equal to 1.

- 1 32. A process for obtaining genotype data/sample allele frequency data as in claim 31, wherein the
2 covering markers are substantially nonevenly distributed across a chromosome or a chromosomal
3 segment.
- 4 33. A process for obtaining genotype data/sample allele frequency data as in claim 31, wherein the
5 covering markers are substantially evenly distributed across a chromosome or a chromosomal
6 segment, and wherein the least common allele frequency of one or more markers is less than 0.4.
- 7 34. A process for obtaining genotype data/sample allele frequency data as in claim 31, wherein the
8 covering markers are substantially evenly distributed across a chromosome or a chromosomal
9 segment; and wherein there is a subgroup of one or more of the covering markers, and each of the
10 markers in the subgroup is chosen without substantial preference for the least common allele frequency
11 of each of the markers in the subgroup being close to 0.5.
- 12 35. A process for obtaining genotype data/sample allele frequency data as in claim 31, wherein the
13 covering markers are substantially evenly distributed across a chromosome or a chromosomal
14 segment, wherein (1) the average chromosomal intermarker distance of the covering markers is greater
15 than 2 cM or the equivalent thereof and the least common allele frequency of one or more of the
16 covering markers is less than 0.3 or wherein (2) the least common allele frequency of one or more of
17 the covering markers is less than 0.2.
- 18 36. A process for obtaining genotype data/sample allele frequency data as in claim 31, wherein the
19 covering markers are substantially evenly distributed across a chromosome or a chromosomal
20 segment, wherein the average chromosomal intermarker distance of the covering markers is less than
21 or equal to 2 cM or the equivalent thereof, and wherein the conditional probability the covering markers
22 were chosen essentially randomly from substantially the known set of bi-allelic markers with least
23 common allele frequencies between 0.2 inclusive and 0.5 inclusive is less than about 10 percent;
24 wherein the conditional probability is substantially conditional on (1) the approximate chromosomal
25 distribution of the covering markers, (2) the marker type of each covering marker and (3) there being N
26 or more covering markers in each cell of the matrix.
- 27 37. A process for obtaining genotype data/sample allele frequency data as in claim 31, wherein the
28 covering markers are substantially evenly distributed across a chromosome or a chromosomal
29 segment, wherein the average chromosomal intermarker distance of the covering markers is greater
30 than 2 cM or the equivalent thereof; and wherein the conditional probability the covering markers were
31 chosen essentially randomly from substantially the known set of bi-allelic markers with least common
32 allele frequencies between 0.3 inclusive and 0.5 inclusive is less than about 10 percent; wherein the
33 conditional probability is substantially conditional on (1) the approximate chromosomal distribution of
34 the covering markers, (2) the marker type of each covering marker and (3) there being N or more
35 covering markers in each cell of the matrix.
- 36 38. A process for obtaining genotype data/sample allele frequency data as in claim 36, wherein the
37 chromosome or the chromosomal segment consists essentially of a set of nonoverlapping chromosome
38 segments of substantially equal length, and wherein one and only one covering marker is located on
39 each of 80 percent or more of the chromosome segments of the set, and wherein zero or two and only
40 two covering markers are located on each of 20 percent or less of the chromosome segments of the

1 set, and wherein each chromosome segment with zero covering markers located thereon is bordered
2 only by chromosome segments with one or two covering markers located thereon, and wherein each
3 chromosome segment with two covering markers located thereon is bordered only by chromosome
4 segments with one or zero covering markers located thereon; and wherein the conditional probability
5 the covering markers were chosen essentially randomly from substantially the known set of bi-allelic
6 markers with least common allele frequencies between 0.2 inclusive and 0.5 inclusive is less than about
7 10 percent; wherein the conditional probability is substantially conditional on (1) the chromosomal
8 distribution of the covering markers on the chromosome segments of the set, (2) the marker type of
9 each covering marker and (3) there being N or more covering markers in each cell of the matrix.

10 39. A process for obtaining genotype data/sample allele frequency data as in claim 37, wherein the
11 chromosome or the chromosomal segment consists essentially of a set of nonoverlapping chromosome
12 segments of substantially equal length, and wherein one and only one covering marker is located on
13 each of 80 percent or more of the chromosome segments of the set, and wherein zero or two and only
14 two covering markers are located on each of 20 percent or less of the chromosome segments of the
15 set, and wherein each chromosome segment with zero covering markers located thereon is bordered
16 only by chromosome segments with one or two covering markers located thereon, and wherein each
17 chromosome segment with two covering markers located thereon is bordered only by chromosome
18 segments with one or zero covering markers located thereon; and wherein the conditional probability
19 the covering markers were chosen essentially randomly from substantially the known set of bi-allelic
20 markers with least common allele frequencies between 0.3 inclusive and 0.5 inclusive is less than about
21 10 percent; wherein the conditional probability is substantially conditional on (1) the chromosomal
22 distribution of the covering markers on the chromosome segments of the set, (2) the marker type of
23 each covering marker and (3) there being N or more covering markers in each cell of the matrix.

24 40. A process for obtaining genotype data/sample allele frequency data as in claim 31, wherein the
25 covering markers are substantially evenly distributed across a chromosome or a chromosomal
26 segment, wherein the average chromosomal intermarker distance of the covering markers is less than
27 or equal to 2 cM or the equivalent thereof, and wherein collection C is essentially the collection of
28 known groups of bi-allelic markers with least common allele frequencies between 0.2 inclusive and 0.5
29 inclusive that are substantially similar to the covering markers as a group; wherein a group of bi-allelic
30 markers is a member of collection C if and only if the group substantially meets criteria (1), (2), (3) and
31 (4): (1) each marker in the group is chosen from substantially the known set of bi-allelic markers with
32 least common allele frequencies between 0.2 inclusive and 0.5 inclusive, (2) the number of markers in
33 the group is the same as the number of covering markers, (3) the chromosomal distribution of the group
34 of markers and the covering markers is substantially similar, and (4) the marker type of each group
35 marker and the covering marker with substantially the same chromosomal location is the same; wherein
36 a group that is a member of collection C substantially meets criterion (5) if and only if (5) there are N or
37 more of the group markers in each cell of the matrix; wherein P is essentially the proportion of groups in
38 collection C that meet criterion (5); wherein P is less than about 90 percent.

39 41. A process for obtaining genotype data/sample allele frequency data as in claim 31, wherein the
40 covering markers are substantially evenly distributed across a chromosome or a chromosomal

1 segment, wherein the average chromosomal intermarker distance of the covering markers is greater
2 than 2 cM or the equivalent thereof, and wherein collection C is essentially the collection of known
3 groups of bi-allelic markers with least common allele frequencies between 0.3 inclusive and 0.5
4 inclusive that are substantially similar to the covering markers as a group; wherein a group of bi-allelic
5 markers is a member of collection C if and only if the group substantially meets criteria (1), (2), (3) and
6 (4): (1) each marker in the group is chosen from substantially the known set of bi-allelic markers with
7 least common allele frequencies between 0.3 inclusive and 0.5 inclusive, (2) the number of markers in
8 the group is the same as the number of covering markers, (3) the chromosomal distribution of the group
9 of markers and the covering markers is substantially similar, and (4) the marker type of each group
10 marker and the covering marker with substantially the same chromosomal location is the same; wherein
11 a group that is a member of collection C substantially meets criterion (5) if and only if (5) there are N or
12 more of the group markers in each cell of the matrix; wherein P is essentially the proportion of groups in
13 collection C that meet criterion (5); wherein P is less than about 90 percent.

14 42. A process for obtaining genotype data/sample allele frequency data as in claim 31, wherein the
15 covering markers are substantially evenly distributed across a chromosome or a chromosomal
16 segment, wherein the average chromosomal intermarker distance of the covering markers is less than
17 or equal to 2 cM or the equivalent thereof; wherein the chromosome or the chromosomal segment
18 consists essentially of a set of nonoverlapping chromosome segments of substantially equal length, and
19 wherein one and only one covering marker is located on each of 80 percent or more of the
20 chromosome segments of the set, and wherein zero or two and only two covering markers are located
21 on each of 20 percent or less of the chromosome segments of the set, and wherein each chromosome
22 segment with zero covering markers located thereon is bordered only by chromosome segments with
23 one or two covering markers located thereon, and wherein each chromosome segment with two
24 covering markers located thereon is bordered only by chromosome segments with one or zero covering
25 markers located thereon; wherein collection D is essentially the collection of known groups of bi-allelic
26 markers with least common allele frequencies between 0.2 inclusive and 0.5 inclusive that are
27 substantially similar to the covering markers as a group; wherein a group of bi-allelic markers is a
28 member of collection D if and only if the group substantially meets criteria (1), (2), and (3): (1) each
29 marker in the group is chosen from substantially the known set of bi-allelic markers with least common
30 allele frequencies between 0.2 inclusive and 0.5 inclusive, (2) the number of covering markers and the
31 number of group markers located on each chromosome segment of the set is the same, and (3) there is
32 a group marker of the same type as each covering marker located on the same chromosome segment
33 of the set as each covering marker; wherein a group that is a member of collection D substantially
34 meets criterion (5) if and only if (5) there are N or more of the group markers in each cell of the matrix;
35 wherein P is essentially the proportion of groups in collection D that meet criterion (5); wherein P is less
36 than about 90 percent.

37 43. A process for obtaining genotype data/sample allele frequency data as in claim 31, wherein the
38 covering markers are substantially evenly distributed across a chromosome or a chromosomal
39 segment, wherein the average chromosomal intermarker distance of the covering markers is greater
40 than 2 cM or the equivalent thereof; wherein the chromosome or the chromosomal segment consists

1 essentially of a set of nonoverlapping chromosome segments of substantially equal length, and wherein
2 one and only one covering marker is located on each of 80 percent or more of the chromosome
3 segments of the set, and wherein zero or two and only two covering markers are located on each of 20
4 percent or less of the chromosome segments of the set, and wherein each chromosome segment with
5 zero covering markers located thereon is bordered only by chromosome segments with one or two
6 covering markers located thereon, and wherein each chromosome segment with two covering markers
7 located thereon is bordered only by chromosome segments with one or zero covering markers located
8 thereon; wherein collection D is essentially the collection of known groups of bi-allelic markers with
9 least common allele frequencies between 0.3 inclusive and 0.5 inclusive that are substantially similar to
10 the covering markers as a group; wherein a group of bi-allelic markers is a member of collection D if
11 and only if the group substantially meets criteria (1), (2), and (3): (1) each marker in the group is chosen
12 from substantially the known set of bi-allelic markers with least common allele frequencies between 0.3
13 inclusive and 0.5 inclusive, (2) the number of covering markers and the number of group markers
14 located on each chromosome segment of the set is the same, and (3) there is a group marker of the
15 same type as each covering marker located on the same chromosome segment of the set as each
16 covering marker; wherein a group that is a member of collection D substantially meets criterion (5) if
17 and only if (5) there are N or more of the group markers in each cell of the matrix; wherein P is
18 essentially the proportion of groups in collection D that meet criterion (5); wherein P is less than about
19 90 percent.

20
21 44. A process for obtaining genotype data/sample allele frequency data as in claim 30, wherein δ is less
22 than or equal to about [1 cM, 0.15] or the equivalent thereof.

23 45. A process for obtaining genotype data/sample allele frequency data as in claim 30, wherein (1) the
24 covering markers are substantially nonevenly distributed across a chromosome or a chromosomal
25 segment or (2) wherein the covering markers are substantially evenly distributed across a chromosome
26 or a chromosomal segment, and wherein the least common allele frequency of one or more markers is
27 less than 0.4 or (3) wherein the covering markers are substantially evenly distributed across a
28 chromosome or a chromosomal segment; and wherein there is a subgroup of one or more of the
29 covering markers, and each of the markers in the subgroup is chosen without substantial preference for
30 the least common allele frequency of each of the markers in the subgroup being close to 0.5.

31 46. A process for obtaining genotype data/sample allele frequency data as in any one of claims 29-45,
32 wherein there is a group of covering markers, and the markers in the group are a majority of the
33 covering markers, and each marker in the group is an SNP, or a bi-allelic marker equivalent formed
34 only from one or more SNPs.

35 47. A process for obtaining genotype data/sample allele frequency data as in claim 31 wherein L_{MC} is
36 less than or equal to about 250,000 bp or the equivalent thereof, W_{MC} is less than or equal to about
37 0.15, wherein the species is human being, wherein the same statistical linkage test based on allelic
38 association is chosen for each covering marker in step b) and wherein there is a group of covering
39 markers, and the markers in the group are a majority of the covering markers, and each marker in the
40 group is an SNP, or a bi-allelic marker equivalent formed only from one or more SNPs.

1 48. An apparatus for obtaining genotype data/sample allele frequency data for each bi-allelic marker of
2 a group of two or more bi-allelic covering markers in the chromosomal DNA of one or more individuals
3 of a sample, each individual in the sample being a member of the same species, comprising: means to
4 practice each of the steps of a process as in any one of the claims 29-47.

5 49. A process for obtaining genotype data/sample allele frequency data as in any one of claims 29-47,
6 wherein the process comprises a computer program.

7 50. An apparatus as in claim 48, wherein the apparatus comprises a computer, the computer being
8 supplied with proper data and instructions.

9
10 51. The use of one or more copies of a set of oligonucleotides to determine genotype data/sample
11 allele frequency data for each bi-allelic marker of a group of two or more bi-allelic covering markers for
12 one or more individuals, each individual being a member of the same species, wherein the group of
13 covering markers systematically cover a CL-F region, the CL-F region being a collection of points on a
14 two-dimensional plane, the two-dimensional plane having the two orthogonal dimensions of
15 chromosomal location and least common allele frequency.

16 52. The use as in claim 51, wherein the CL-F region is N covered to within a CL-F distance δ by the two
17 or more bi-allelic covering markers, so that each point in the region is within the CL-F distance δ of N or
18 more of the covering markers, wherein δ is equal to about $[\delta_{CL}, 0.25]$ or the equivalent thereof, δ_{CL} is
19 equal to the largest chromosomal length, computed by any method, for which linkage disequilibrium has
20 been observed between any polymorphisms in any population of the species, N is an integer greater
21 than or equal to 1.

22 53. The use as in claim 51, wherein the CL-F region comprises a CL-F matrix, the sum of the number of
23 columns and rows in the matrix being greater than or equal to three, each cell of the matrix being of
24 length L_{MC} and width W_{MC} , and L_{MC} being less than or equal to about δ_{CL} , and W_{MC} being less than or
25 equal to about 0.25, δ_{CL} is equal to the largest chromosomal length, computed by any method, for which
26 linkage disequilibrium has been observed between any polymorphisms in any population of the species,
27 there being N or more covering markers in each cell of the matrix and N is an integer greater than or
28 equal to 1.

29 54. The use as in claim 53, wherein the covering markers are substantially nonevenly distributed across
30 a chromosome or a chromosomal segment.

31 55. The use as in claim 53, wherein the covering markers are substantially evenly distributed across a
32 chromosome or a chromosomal segment, and wherein the least common allele frequency of one or
33 more markers is less than 0.4.

34 56. The use as in claim 53, wherein the covering markers are substantially evenly distributed across a
35 chromosome or a chromosomal segment; and wherein there is a subgroup of one or more of the
36 covering markers, and each of the markers in the subgroup is chosen without substantial preference for
37 the least common allele frequency of each of the markers in the subgroup being close to 0.5.

38 57. The use as in claim 53, wherein the covering markers are substantially evenly distributed across a
39 chromosome or a chromosomal segment, wherein (1) the average chromosomal intermarker distance
40 of the covering markers is greater than 2 cM or the equivalent thereof and the least common allele

1 frequency of one or more of the covering markers is less than 0.3 or wherein (2) the least common
2 allele frequency of one or more of the covering markers is less than 0.2.

3 58. The use as in claim 53, wherein the covering markers are substantially evenly distributed across a
4 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
5 the covering markers is less than or equal to 2 cM or the equivalent thereof, and wherein the conditional
6 probability the covering markers were chosen essentially randomly from substantially the known set of
7 bi-allelic markers with least common allele frequencies between 0.2 inclusive and 0.5 inclusive is less
8 than about 10 percent; wherein the conditional probability is substantially conditional on (1) the
9 approximate chromosomal distribution of the covering markers, (2) the marker type of each covering
10 marker and (3) there being N or more covering markers in each cell of the matrix.

11 59. The use as in claim 53, wherein the covering markers are substantially evenly distributed across a
12 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
13 the covering markers is greater than 2 cM or the equivalent thereof; and wherein the conditional
14 probability the covering markers were chosen essentially randomly from substantially the known set of
15 bi-allelic markers with least common allele frequencies between 0.3 inclusive and 0.5 inclusive is less
16 than about 10 percent; wherein the conditional probability is substantially conditional on (1) the
17 approximate chromosomal distribution of the covering markers, (2) the marker type of each covering
18 marker and (3) there being N or more covering markers in each cell of the matrix.

19 60. The use as in claim 58, wherein the chromosome or the chromosomal segment consists essentially
20 of a set of nonoverlapping chromosome segments of substantially equal length, and wherein one and
21 only one covering marker is located on each of 80 percent or more of the chromosome segments of the
22 set, and wherein zero or two and only two covering markers are located on each of 20 percent or less
23 of the chromosome segments of the set, and wherein each chromosome segment with zero covering
24 markers located thereon is bordered only by chromosome segments with one or two covering markers
25 located thereon, and wherein each chromosome segment with two covering markers located thereon is
26 bordered only by chromosome segments with one or zero covering markers located thereon; and
27 wherein the conditional probability the covering markers were chosen essentially randomly from
28 substantially the known set of bi-allelic markers with least common allele frequencies between 0.2
29 inclusive and 0.5 inclusive is less than about 10 percent; wherein the conditional probability is
30 substantially conditional on (1) the chromosomal distribution of the covering markers on the
31 chromosome segments of the set, (2) the marker type of each covering marker and (3) there being N or
32 more covering markers in each cell of the matrix.

33 61. The use as in claim 59, wherein the chromosome or the chromosomal segment consists essentially
34 of a set of nonoverlapping chromosome segments of substantially equal length, and wherein one and
35 only one covering marker is located on each of 80 percent or more of the chromosome segments of the
36 set, and wherein zero or two and only two covering markers are located on each of 20 percent or less
37 of the chromosome segments of the set, and wherein each chromosome segment with zero covering
38 markers located thereon is bordered only by chromosome segments with one or two covering markers
39 located thereon, and wherein each chromosome segment with two covering markers located thereon is
40 bordered only by chromosome segments with one or zero covering markers located thereon; and

1 wherein the conditional probability the covering markers were chosen essentially randomly from
2 substantially the known set of bi-allelic markers with least common allele frequencies between 0.3
3 inclusive and 0.5 inclusive is less than about 10 percent; wherein the conditional probability is
4 substantially conditional on (1) the chromosomal distribution of the covering markers on the
5 chromosome segments of the set, (2) the marker type of each covering marker and (3) there being N or
6 more covering markers in each cell of the matrix.

7 62. The use as in claim 53, wherein the covering markers are substantially evenly distributed across a
8 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
9 the covering markers is less than or equal to 2 cM or the equivalent thereof, and wherein collection C is
10 essentially the collection of known groups of bi-allelic markers with least common allele frequencies
11 between 0.2 inclusive and 0.5 inclusive that are substantially similar to the covering markers as a
12 group; wherein a group of bi-allelic markers is a member of collection C if and only if the group
13 substantially meets criteria (1), (2), (3) and (4): (1) each marker in the group is chosen from
14 substantially the known set of bi-allelic markers with least common allele frequencies between 0.2
15 inclusive and 0.5 inclusive, (2) the number of markers in the group is the same as the number of
16 covering markers, (3) the chromosomal distribution of the group of markers and the covering markers is
17 substantially similar, and (4) the marker type of each group marker and the covering marker with
18 substantially the same chromosomal location is the same; wherein a group that is a member of
19 collection C substantially meets criterion (5) if and only if (5) there are N or more of the group markers
20 in each cell of the matrix; wherein P is essentially the proportion of groups in collection C that meet
21 criterion (5); wherein P is less than about 90 percent.

22 63. The use as in claim 53, wherein the covering markers are substantially evenly distributed across a
23 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
24 the covering markers is greater than 2 cM or the equivalent thereof, and wherein collection C is
25 essentially the collection of known groups of bi-allelic markers with least common allele frequencies
26 between 0.3 inclusive and 0.5 inclusive that are substantially similar to the covering markers as a
27 group; wherein a group of bi-allelic markers is a member of collection C if and only if the group
28 substantially meets criteria (1), (2), (3) and (4): (1) each marker in the group is chosen from
29 substantially the known set of bi-allelic markers with least common allele frequencies between 0.3
30 inclusive and 0.5 inclusive, (2) the number of markers in the group is the same as the number of
31 covering markers, (3) the chromosomal distribution of the group of markers and the covering markers is
32 substantially similar, and (4) the marker type of each group marker and the covering marker with
33 substantially the same chromosomal location is the same; wherein a group that is a member of
34 collection C substantially meets criterion (5) if and only if (5) there are N or more of the group markers
35 in each cell of the matrix; wherein P is essentially the proportion of groups in collection C that meet
36 criterion (5); wherein P is less than about 90 percent.

37 64. The use as in claim 53, wherein the covering markers are substantially evenly distributed across a
38 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
39 the covering markers is less than or equal to 2 cM or the equivalent thereof; wherein the chromosome
40 or the chromosomal segment consists essentially of a set of nonoverlapping chromosome segments of

1 substantially equal length, and wherein one and only one covering marker is located on each of 80
2 percent or more of the chromosome segments of the set, and wherein zero or two and only two
3 covering markers are located on each of 20 percent or less of the chromosome segments of the set,
4 and wherein each chromosome segment with zero covering markers located thereon is bordered only
5 by chromosome segments with one or two covering markers located thereon, and wherein each
6 chromosome segment with two covering markers located thereon is bordered only by chromosome
7 segments with one or zero covering markers located thereon; wherein collection D is essentially the
8 collection of known groups of bi-allelic markers with least common allele frequencies between 0.2
9 inclusive and 0.5 inclusive that are substantially similar to the covering markers as a group; wherein a
10 group of bi-allelic markers is a member of collection D if and only if the group substantially meets
11 criteria (1), (2), and (3): (1) each marker in the group is chosen from substantially the known set of bi-
12 allelic markers with least common allele frequencies between 0.2 inclusive and 0.5 inclusive, (2) the
13 number of covering markers and the number of group markers located on each chromosome segment
14 of the set is the same, and (3) there is a group marker of the same type as each covering marker
15 located on the same chromosome segment of the set as each covering marker; wherein a group that is
16 a member of collection D substantially meets criterion (5) if and only if (5) there are N or more of the
17 group markers in each cell of the matrix; wherein P is essentially the proportion of groups in collection D
18 that meet criterion (5); wherein P is less than about 90 percent.

19 65. The use as in claim 53, wherein the covering markers are substantially evenly distributed across a
20 chromosome or a chromosomal segment, wherein the average chromosomal intermarker distance of
21 the covering markers is greater than 2 cM or the equivalent thereof; wherein the chromosome or the
22 chromosomal segment consists essentially of a set of nonoverlapping chromosome segments of
23 substantially equal length, and wherein one and only one covering marker is located on each of 80
24 percent or more of the chromosome segments of the set, and wherein zero or two and only two
25 covering markers are located on each of 20 percent or less of the chromosome segments of the set,
26 and wherein each chromosome segment with zero covering markers located thereon is bordered only
27 by chromosome segments with one or two covering markers located thereon, and wherein each
28 chromosome segment with two covering markers located thereon is bordered only by chromosome
29 segments with one or zero covering markers located thereon; wherein collection D is essentially the
30 collection of known groups of bi-allelic markers with least common allele frequencies between 0.3
31 inclusive and 0.5 inclusive that are substantially similar to the covering markers as a group; wherein a
32 group of bi-allelic markers is a member of collection D if and only if the group substantially meets
33 criteria (1), (2), and (3): (1) each marker in the group is chosen from substantially the known set of bi-
34 allelic markers with least common allele frequencies between 0.3 inclusive and 0.5 inclusive, (2) the
35 number of covering markers and the number of group markers located on each chromosome segment
36 of the set is the same, and (3) there is a group marker of the same type as each covering marker
37 located on the same chromosome segment of the set as each covering marker; wherein a group that is
38 a member of collection D substantially meets criterion (5) if and only if (5) there are N or more of the
39 group markers in each cell of the matrix; wherein P is essentially the proportion of groups in collection D
40 that meet criterion (5); wherein P is less than about 90 percent.

1 66. The use as in claim 52, wherein δ is less than or equal to about [1 cM, 0.15] or the equivalent
2 thereof.

3 67. The use as in claim 52, wherein (1) the covering markers are substantially nonevenly distributed
4 across a chromosome or a chromosomal segment or (2) wherein the covering markers are substantially
5 evenly distributed across a chromosome or a chromosomal segment, and wherein the least common
6 allele frequency of one or more markers is less than 0.4 or (3) wherein the covering markers are
7 substantially evenly distributed across a chromosome or a chromosomal segment; and wherein there is
8 a subgroup of one or more of the covering markers, and each of the markers in the subgroup is chosen
9 without substantial preference for the least common allele frequency of each of the markers in the
10 subgroup being close to 0.5.

11 68. The use as in any one of claims 51-67, wherein there is a group of covering markers, and the
12 markers in the group are a majority of the covering markers, and each marker in the group is an SNP,
13 or a bi-allelic marker equivalent formed only from one or more SNPs.

14 69. The use as in claim 53 wherein L_{MC} is less than or equal to about 250,000 bp or the equivalent
15 thereof, W_{MC} is less than or equal to about 0.15, wherein the species is human being, wherein the same
16 statistical linkage test based on allelic association is chosen for each covering marker in step b) and
17 wherein there is a group of covering markers, and the markers in the group are a majority of the
18 covering markers, and each marker in the group is an SNP, or a bi-allelic marker equivalent formed
19 only from one or more SNPs.

20
21 70. One or more copies of a set of oligonucleotides, the set of oligonucleotides being complementary to
22 a group of two or more bi-allelic covering markers of the same species, wherein the group of covering
23 markers systematically cover a CL-F region, the CL-F region being a collection of points on a two-
24 dimensional plane, the two-dimensional plane having the two orthogonal dimensions of chromosomal
25 location and least common allele frequency.

26 71. One or more copies of a set of oligonucleotides as in claim 70, wherein the CL-F region is N
27 covered to within a CL-F distance δ by the two or more bi-allelic covering markers, so that each point in
28 the region is within the CL-F distance δ of N or more of the covering markers, wherein δ is equal to
29 about [δ_{CL} , 0.25] or the equivalent thereof, δ_{CL} is equal to the largest chromosomal length, computed by
30 any method, for which linkage disequilibrium has been observed between any polymorphisms in any
31 population of the species, N is an integer greater than or equal to 1.

32 72. One or more copies of a set of oligonucleotides as in claim 70, wherein the CL-F region comprises
33 a CL-F matrix, the sum of the number of columns and rows in the matrix being greater than or equal to
34 three, each cell of the matrix being of length L_{MC} and width W_{MC} , and L_{MC} being less than or equal to
35 about δ_{CL} , and W_{MC} being less than or equal to about 0.25, δ_{CL} is equal to the largest chromosomal
36 length, computed by any method, for which linkage disequilibrium has been observed between any
37 polymorphisms in any population of the species, there being N or more covering markers in each cell of
38 the matrix and N is an integer greater than or equal to 1.

39 73. One or more copies of a set of oligonucleotides as in claim 72, wherein the covering markers are
40 substantially nonevenly distributed across a chromosome or a chromosomal segment.

1 74. One or more copies of a set of oligonucleotides as in claim 72, wherein the covering markers are
2 substantially evenly distributed across a chromosome or a chromosomal segment, and wherein the
3 least common allele frequency of one or more markers is less than 0.4.

4 75. One or more copies of a set of oligonucleotides as in claim 72, wherein the covering markers are
5 substantially evenly distributed across a chromosome or a chromosomal segment; and wherein there is
6 a subgroup of one or more of the covering markers, and each of the markers in the subgroup is chosen
7 without substantial preference for the least common allele frequency of each of the markers in the
8 subgroup being close to 0.5.

9 76. One or more copies of a set of oligonucleotides as in claim 72, wherein the covering markers are
10 substantially evenly distributed across a chromosome or a chromosomal segment, wherein (1) the
11 average chromosomal intermarker distance of the covering markers is greater than 2 cM or the
12 equivalent thereof and the least common allele frequency of one or more of the covering markers is
13 less than 0.3 or wherein (2) the least common allele frequency of one or more of the covering markers
14 is less than 0.2.

15 77. One or more copies of a set of oligonucleotides as in claim 72, wherein the covering markers are
16 substantially evenly distributed across a chromosome or a chromosomal segment, wherein the average
17 chromosomal intermarker distance of the covering markers is less than or equal to 2 cM or the
18 equivalent thereof, and wherein the conditional probability the covering markers were chosen
19 essentially randomly from substantially the known set of bi-allelic markers with least common allele
20 frequencies between 0.2 inclusive and 0.5 inclusive is less than about 10 percent; wherein the
21 conditional probability is substantially conditional on (1) the approximate chromosomal distribution of
22 the covering markers, (2) the marker type of each covering marker and (3) there being N or more
23 covering markers in each cell of the matrix.

24 78. One or more copies of a set of oligonucleotides as in claim 72, wherein the covering markers are
25 substantially evenly distributed across a chromosome or a chromosomal segment, wherein the average
26 chromosomal intermarker distance of the covering markers is greater than 2 cM or the equivalent
27 thereof; and wherein the conditional probability the covering markers were chosen essentially randomly
28 from substantially the known set of bi-allelic markers with least common allele frequencies between 0.3
29 inclusive and 0.5 inclusive is less than about 10 percent; wherein the conditional probability is
30 substantially conditional on (1) the approximate chromosomal distribution of the covering markers, (2)
31 the marker type of each covering marker and (3) there being N or more covering markers in each cell of
32 the matrix.

33 79. One or more copies of a set of oligonucleotides as in claim 77, wherein the chromosome or the
34 chromosomal segment consists essentially of a set of nonoverlapping chromosome segments of
35 substantially equal length, and wherein one and only one covering marker is located on each of 80
36 percent or more of the chromosome segments of the set, and wherein zero or two and only two
37 covering markers are located on each of 20 percent or less of the chromosome segments of the set,
38 and wherein each chromosome segment with zero covering markers located thereon is bordered only
39 by chromosome segments with one or two covering markers located thereon, and wherein each
40 chromosome segment with two covering markers located thereon is bordered only by chromosome

1 segments with one or zero covering markers located thereon; and wherein the conditional probability
2 the covering markers were chosen essentially randomly from substantially the known set of bi-allelic
3 markers with least common allele frequencies between 0.2 inclusive and 0.5 inclusive is less than about
4 10 percent; wherein the conditional probability is substantially conditional on (1) the chromosomal
5 distribution of the covering markers on the chromosome segments of the set, (2) the marker type of
6 each covering marker and (3) there being N or more covering markers in each cell of the matrix.
7 80. One or more copies of a set of oligonucleotides as in claim 78, wherein the chromosome or the
8 chromosomal segment consists essentially of a set of nonoverlapping chromosome segments of
9 substantially equal length, and wherein one and only one covering marker is located on each of 80
10 percent or more of the chromosome segments of the set, and wherein zero or two and only two
11 covering markers are located on each of 20 percent or less of the chromosome segments of the set,
12 and wherein each chromosome segment with zero covering markers located thereon is bordered only
13 by chromosome segments with one or two covering markers located thereon, and wherein each
14 chromosome segment with two covering markers located thereon is bordered only by chromosome
15 segments with one or zero covering markers located thereon; and wherein the conditional probability
16 the covering markers were chosen essentially randomly from substantially the known set of bi-allelic
17 markers with least common allele frequencies between 0.3 inclusive and 0.5 inclusive is less than about
18 10 percent; wherein the conditional probability is substantially conditional on (1) the chromosomal
19 distribution of the covering markers on the chromosome segments of the set, (2) the marker type of
20 each covering marker and (3) there being N or more covering markers in each cell of the matrix.
21 81. One or more copies of a set of oligonucleotides as in claim 72, wherein the covering markers are
22 substantially evenly distributed across a chromosome or a chromosomal segment, wherein the average
23 chromosomal intermarker distance of the covering markers is less than or equal to 2 cM or the
24 equivalent thereof, and wherein collection C is essentially the collection of known groups of bi-allelic
25 markers with least common allele frequencies between 0.2 inclusive and 0.5 inclusive that are
26 substantially similar to the covering markers as a group; wherein a group of bi-allelic markers is a
27 member of collection C if and only if the group substantially meets criteria (1), (2), (3) and (4): (1) each
28 marker in the group is chosen from substantially the known set of bi-allelic markers with least common
29 allele frequencies between 0.2 inclusive and 0.5 inclusive, (2) the number of markers in the group is the
30 same as the number of covering markers, (3) the chromosomal distribution of the group of markers and
31 the covering markers is substantially similar, and (4) the marker type of each group marker and the
32 covering marker with substantially the same chromosomal location is the same; wherein a group that is
33 a member of collection C substantially meets criterion (5) if and only if (5) there are N or more of the
34 group markers in each cell of the matrix; wherein P is essentially the proportion of groups in collection C
35 that meet criterion (5); wherein P is less than about 90 percent.
36 82. One or more copies of a set of oligonucleotides as in claim 72, wherein the covering markers are
37 substantially evenly distributed across a chromosome or a chromosomal segment, wherein the average
38 chromosomal intermarker distance of the covering markers is greater than 2 cM or the equivalent
39 thereof, and wherein collection C is essentially the collection of known groups of bi-allelic markers with
40 least common allele frequencies between 0.3 inclusive and 0.5 inclusive that are substantially similar to

1 the covering markers as a group; wherein a group of bi-allelic markers is a member of collection C if
2 and only if the group substantially meets criteria (1), (2), (3) and (4): (1) each marker in the group is
3 chosen from substantially the known set of bi-allelic markers with least common allele frequencies
4 between 0.3 inclusive and 0.5 inclusive, (2) the number of markers in the group is the same as the
5 number of covering markers, (3) the chromosomal distribution of the group of markers and the covering
6 markers is substantially similar, and (4) the marker type of each group marker and the covering marker
7 with substantially the same chromosomal location is the same; wherein a group that is a member of
8 collection C substantially meets criterion (5) if and only if (5) there are N or more of the group markers
9 in each cell of the matrix; wherein P is essentially the proportion of groups in collection C that meet
10 criterion (5); wherein P is less than about 90 percent.

11 83. One or more copies of a set of oligonucleotides as in claim 72, wherein the covering markers are
12 substantially evenly distributed across a chromosome or a chromosomal segment, wherein the average
13 chromosomal intermarker distance of the covering markers is less than or equal to 2 cM or the
14 equivalent thereof; wherein the chromosome or the chromosomal segment consists essentially of a set
15 of nonoverlapping chromosome segments of substantially equal length, and wherein one and only one
16 covering marker is located on each of 80 percent or more of the chromosome segments of the set, and
17 wherein zero or two and only two covering markers are located on each of 20 percent or less of the
18 chromosome segments of the set, and wherein each chromosome segment with zero covering markers
19 located thereon is bordered only by chromosome segments with one or two covering markers located
20 thereon, and wherein each chromosome segment with two covering markers located thereon is
21 bordered only by chromosome segments with one or zero covering markers located thereon; wherein
22 collection D is essentially the collection of known groups of bi-allelic markers with least common allele
23 frequencies between 0.2 inclusive and 0.5 inclusive that are substantially similar to the covering
24 markers as a group; wherein a group of bi-allelic markers is a member of collection D if and only if the
25 group substantially meets criteria (1), (2), and (3): (1) each marker in the group is chosen from
26 substantially the known set of bi-allelic markers with least common allele frequencies between 0.2
27 inclusive and 0.5 inclusive, (2) the number of covering markers and the number of group markers
28 located on each chromosome segment of the set is the same, and (3) there is a group marker of the
29 same type as each covering marker located on the same chromosome segment of the set as each
30 covering marker; wherein a group that is a member of collection D substantially meets criterion (5) if
31 and only if (5) there are N or more of the group markers in each cell of the matrix; wherein P is
32 essentially the proportion of groups in collection D that meet criterion (5); wherein P is less than about
33 90 percent.

34 84. One or more copies of a set of oligonucleotides as in claim 72, wherein the covering markers are
35 substantially evenly distributed across a chromosome or a chromosomal segment, wherein the average
36 chromosomal intermarker distance of the covering markers is greater than 2 cM or the equivalent
37 thereof; wherein the chromosome or the chromosomal segment consists essentially of a set of
38 nonoverlapping chromosome segments of substantially equal length, and wherein one and only one
39 covering marker is located on each of 80 percent or more of the chromosome segments of the set, and
40 wherein zero or two and only two covering markers are located on each of 20 percent or less of the

1 chromosome segments of the set, and wherein each chromosome segment with zero covering markers
2 located thereon is bordered only by chromosome segments with one or two covering markers located
3 thereon, and wherein each chromosome segment with two covering markers located thereon is
4 bordered only by chromosome segments with one or zero covering markers located thereon; wherein
5 collection D is essentially the collection of known groups of bi-allelic markers with least common allele
6 frequencies between 0.3 inclusive and 0.5 inclusive that are substantially similar to the covering
7 markers as a group; wherein a group of bi-allelic markers is a member of collection D if and only if the
8 group substantially meets criteria (1), (2), and (3): (1) each marker in the group is chosen from
9 substantially the known set of bi-allelic markers with least common allele frequencies between 0.3
10 inclusive and 0.5 inclusive, (2) the number of covering markers and the number of group markers
11 located on each chromosome segment of the set is the same, and (3) there is a group marker of the
12 same type as each covering marker located on the same chromosome segment of the set as each
13 covering marker; wherein a group that is a member of collection D substantially meets criterion (5) if
14 and only if (5) there are N or more of the group markers in each cell of the matrix; wherein P is
15 essentially the proportion of groups in collection D that meet criterion (5); wherein P is less than about
16 90 percent.

17 85. One or more copies of a set of oligonucleotides as in claim 71, wherein δ is less than or equal to
18 about [1 cM, 0.15] or the equivalent thereof.

19 86. One or more copies of a set of oligonucleotides as in claim 71, wherein (1) the covering markers are
20 substantially nonevenly distributed across a chromosome or a chromosomal segment or (2) wherein the
21 covering markers are substantially evenly distributed across a chromosome or a chromosomal
22 segment, and wherein the least common allele frequency of one or more markers is less than 0.4 or (3)
23 wherein the covering markers are substantially evenly distributed across a chromosome or a
24 chromosomal segment; and wherein there is a subgroup of one or more of the covering markers, and
25 each of the markers in the subgroup is chosen without substantial preference for the least common
26 allele frequency of each of the markers in the subgroup being close to 0.5.

27 87. One or more copies of a set of oligonucleotides as in any one of claims 70-86, wherein there is a
28 group of covering markers, and the markers in the group are a majority of the covering markers, and
29 each marker in the group is an SNP, or a bi-allelic marker equivalent formed only from one or more
30 SNPs.

31 88. One or more copies of a set of oligonucleotides as in claim 72 wherein L_{MC} is less than or equal to
32 about 250,000 bp or the equivalent thereof, W_{MC} is less than or equal to about 0.15, wherein the
33 species is human being, wherein the same statistical linkage test based on allelic association is chosen
34 for each covering marker in step b) and wherein there is a group of covering markers, and the markers
35 in the group are a majority of the covering markers, and each marker in the group is an SNP, or a bi-
36 allelic marker equivalent formed only from one or more SNPs.

37

Statement under Article 19(1)

Some of the amended claims make use of the phrase "conditional probability", such as claim 11. Some of the amended claims make use of the phrase "proportion of groups", such as claim 14. There are various techniques to calculate or estimate such a probability or such a proportion. These techniques include, but are not necessarily limited to, direct calculation, statistical estimates, and Monte Carlo estimation techniques. Powerful software is available for calculation and statistical estimation for data in matrix format or two-dimensional format. Some such software is available from Cytel Software Corporation, Cambridge, Massachusetts (example: Exact Logistic Regression: Theory and Examples, Mehta CR, Patel NR, Statistics in Medicine, vol 14, 2143-2160(1995). Another example is SAS (SAS Institute Inc., SAS Campus Drive, Cary, North Carolina 27513, USA.; A handbook of statistical analyses using SAS by Brian S. Everitt and Geoff Der, Boca Raton, Fla. : Chapman & Hall/CRC, 1998.). A further example is MATLAB (The MathWorks, Inc. 3 Apple Hill Drive, Natick, Mass. U.S.A. 01760-2098; MATLAB primer by Kermit Sigmon, 4th ed. Boca Raton : CRC Press, c1994.) Statistical techniques include techniques for hypothesis testing, goodness-of-fit and others.

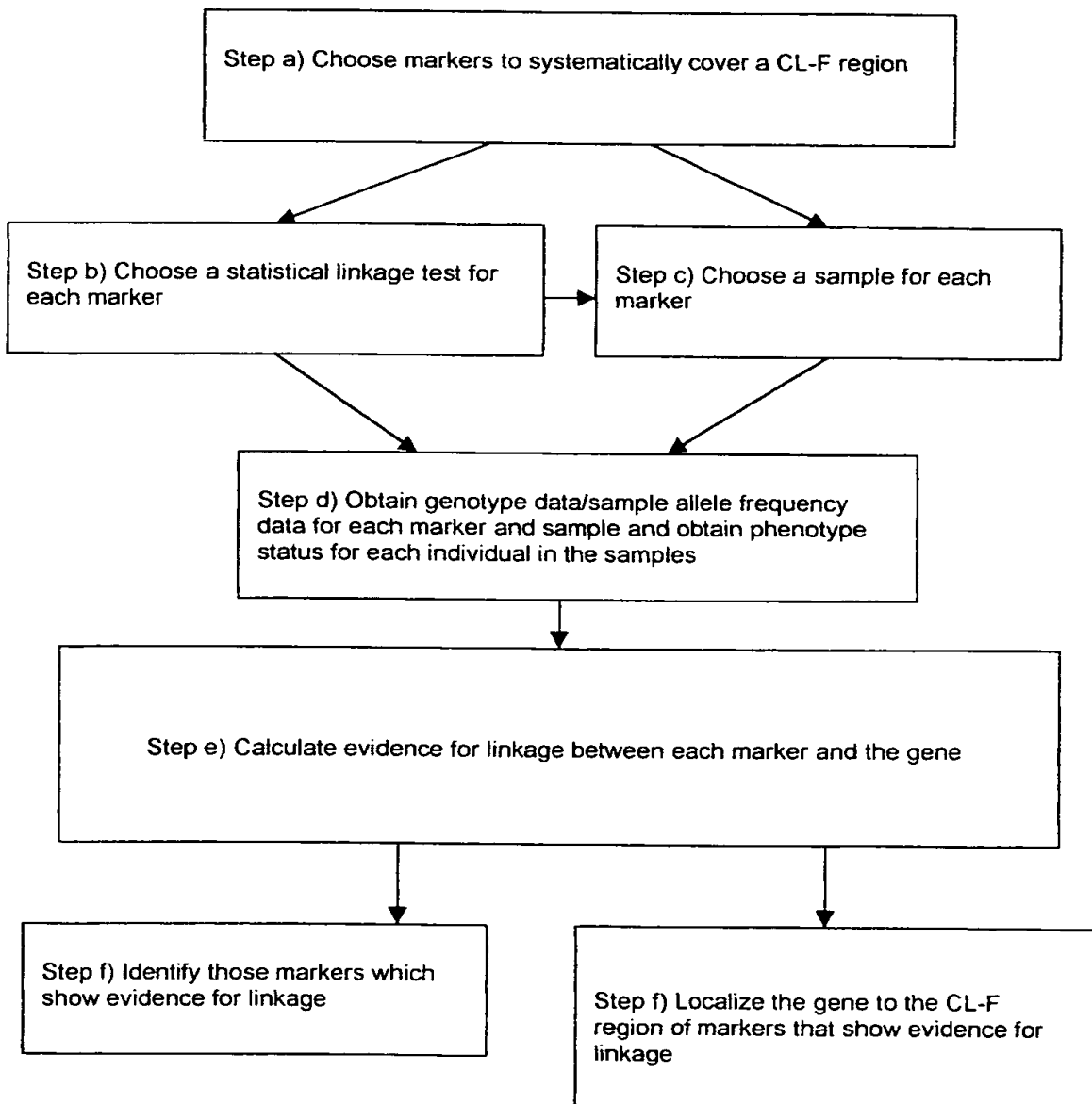
The degree of skill in the art in probability and statistics is great. Indeed the inventor's important equation (Equation 2, page 38) is an equation for P_t , wherein P_t is a binomial probability for parental allele 'transmission' which determines the magnitude of the TDT chi-square statistic. P_s (pages 40-42) is also a binomial probability that determines the magnitude of the ASP test statistic. (see Abstract and Paper: Annals of Human Genetics (1998), 62, 159-179. The abstract is available on the World Wide Web and Internet, including at the journal's website.) Skill in the use of computers in the art is also great (page 25).

Some claims, such as claims 11, 12, 13, 14 and others make use of the phrase "substantially the known set of bi-allelic markers". As pointed out in the description (page 25) information on bi-allelic markers can be gained from sources such as the Whitehead Institute or Marshfield Foundation for Biomedical Research. Similar sources of information on Single Nucleotide Polymorphisms can be obtained from sources given in SNP attack on complex traits, Nature Genetics, volume 20 no. 3, Nov 1998, pp. 217-218.

Some claims, such as claims 11, 12, 13, 14 and others make use of the term "marker type" or similar terminology. As stated in the description, a bi-allelic marker may be an SNP, a microsatellite marker, a bi-allelic marker equivalent formed from one or more true bi-allelic markers. "Marker type" means type of true bi-allelic marker as for example an SNP or a microsatellite; or "marker type" means a bi-allelic marker equivalent of a certain type, such as a bi-allelic marker equivalent formed only from one or more SNPs or a bi-allelic marker equivalent formed only from one or more microsatellites.)

1/1

Drawing #1: Computer Program Flowsheet for Process #1 and Process#1A



SUBSTITUTE SHEET (RULE 26)

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US99/04376

A. CLASSIFICATION OF SUBJECT MATTER

IPC(6) : C12Q 1/68; C07H 21/04

US CL : 435/6, 287.2; 436/94; 536/24.3, 24.31

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

U.S. : 435/6, 287.2; 436/94; 536/24.3, 24.31

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

APS, BIOSIS, CA, MEDLINE

search terms: bi-allelic markers, CL-F, frequency, linkage, matrix, chromosome

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	KRUGLYAK et al. Linkage thresholds for two-stage genome scans. Am. J. Hum. Genet. 1998, Vol. 62, pages 994-996.	2-20, 22-24
A	KRUGLYAK, L. The use of a genetic map of biallelic markers in linkage studies. Nature Genetics. September 1997, Vol. 17, pages 21-24.	2-20, 22-24
A	KRUGLYAK et al. Parametric and nonparametric linkage analysis: A unified multipoint approach. Am. J. Hum. Genet. 1996, Vol. 58, pages 1347-1363.	2-20, 22-24
A	KRUGLYAK et al. Complete multipoint sib-pair analysis of qualitative and quantitative traits. Am. J. Hum. Genet. 1995, Vol. 57, pages 439-454.	2-20, 22-24

☐ Further documents are listed in the continuation of Box C. ☐ See patent family annex.

* Special categories of cited documents.	*T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
A document defining the general state of the art which is not considered to be of particular relevance	*X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
E earlier document published on or after the international filing date	*Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
L document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	*a* document member of the same patent family
O document referring to an oral disclosure, use, exhibition or other means	
P document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search 04 MAY 1999	Date of mailing of the international search report 24 MAY 1999 JOYCE BRINKERS PARALEGAL CH JTB
Name and mailing address of the ISA/US Commissioner of Patents and Trademarks Box PCT Washington, D.C. 20231 Facsimile No. (703) 305-3230	Authorized officer KENNETH R. HORLICK Telephone No. (703) 308-0196

Form PCT/ISA/210 (second sheet)(July 1992)*

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US99/04376

Box I Observations where certain claims were found unsearchable (Continuation of item 1 of first sheet)

This international report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1. ☐ Claims Nos.:
because they relate to subject matter not required to be searched by this Authority, namely:

2. ☒ Claims Nos.: 1
because they relate to parts of the international application that do not comply with the prescribed requirements to such an extent that no meaningful international search can be carried out, specifically:

the claim does not set forth an invention

3. ☒ Claims Nos.: 21 and 25
because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

Box II Observations where unity of invention is lacking (Continuation of item 2 of first sheet)

This International Searching Authority found multiple inventions in this international application, as follows:

1. ☐ As all required additional search fees were timely paid by the applicant, this international search report covers all searchable claims.

2. ☐ As all searchable claims could be searched without effort justifying an additional fee, this Authority did not invite payment of any additional fee.

3. ☐ As only some of the required additional search fees were timely paid by the applicant, this international search report covers only those claims for which fees were paid, specifically claims Nos.:

4. ☐ No required additional search fees were timely paid by the applicant. Consequently, this international search report is restricted to the invention first mentioned in the claims; it is covered by claims Nos.:

Remark on Protest

- ☐ The additional search fees were accompanied by the applicant's protest.
☐ No protest accompanied the payment of additional search fees.

